

Beyond Value-Function Gaps: Improved Instance-Dependent Regret Bounds for Episodic Reinforcement Learning

- Christoph Dann, Teodor Vanislavov Marinov, Mehryar Mohri, Julian Zimmert
- abstract@[open-review](#): We provide improved gap-dependent regret bounds for reinforcement learning in finite episodic Markov decision processes. Compared to prior work, our bounds depend on alternative definitions of gaps. These definitions are based on the insight that, in order to achieve a favorable regret, an algorithm does not need to learn how to behave optimally in states that are not reached by an optimal policy. We prove tighter upper regret bounds for optimistic algorithms and accompany them with new information-theoretic lower bounds for a large class of MDPs. Our results show that optimistic algorithms can not achieve the information-theoretic lower bounds even in deterministic MDPs unless there is a unique optimal policy.

Learning One Representation to Optimize All Rewards

- Ahmed Touati, Yann Ollivier
- abstract@[open-review](#): We introduce the forward-backward (FB) representation of the dynamics of a reward-free Markov decision process. It provides explicit near-optimal policies for any reward specified a posteriori. During an unsupervised phase, we use reward-free interactions with the environment to learn two representations via off-the-shelf deep learning methods and temporal difference (TD) learning. In the test phase, a reward representation is estimated either from reward observations or an explicit reward description (e.g., a target state). The optimal policy for that reward is directly obtained from these representations, with no planning. We assume access to an exploration scheme or replay buffer for the first phase. The corresponding unsupervised loss is well-principled: if training is perfect, the policies obtained are provably optimal for any reward function. With imperfect training, the sub-optimality is proportional to the unsupervised approximation error. The FB representation learns long-range relationships between states and actions, via a predictive occupancy map, without having to synthesize states as in model-based approaches. This is a step towards learning controllable agents in arbitrary black-box stochastic environments. This approach compares well to goal-oriented RL algorithms on discrete and continuous mazes, pixel-based MsPacman, and the FetchReach virtual robot arm. We also illustrate how the agent can immediately adapt to new tasks beyond goal-oriented RL.

Matrix factorisation and the interpretation of geodesic distance

- Nick Whiteley, Annie Gray, Patrick Rubin-Delanchy
- abstract@[open-review](#): Given a graph or similarity matrix, we consider the problem of recovering a notion of true distance between the nodes, and so their true positions. We show that this can be accomplished in two steps: matrix factorisation, followed by nonlinear dimension reduction. This combination is effective because the point cloud obtained in the first step lives close to a manifold in which latent distance is encoded as geodesic distance. Hence, a nonlinear dimension reduction tool, approximating geodesic distance, can recover the latent positions, up to a simple transformation. We give a detailed account of the case where spectral embedding is used, followed by Isomap, and provide encouraging experimental evidence for other combinations of techniques.

UniDoc: Unified Pretraining Framework for Document Understanding

- Jiuxiang Gu, Jason Kuen, Vlad I Morariu, Handong Zhao, Rajiv Jain, Nikolaos Barmpalias, Ani Nenkova, Tong Sun
- abstract@[open-review](#): Document intelligence automates the extraction of information from documents and supports many business applications. Recent self-supervised learning methods on large-scale unlabeled document datasets have opened up promising directions towards reducing annotation efforts by training models with self-supervised objectives. However, most of the existing document pretraining methods are still language-dominated. We present UDoc, a new unified pretraining framework for document understanding. UDoc is designed to support most document understanding tasks, extending the Transformer to take multimodal embeddings as input. Each input element is composed of words and visual features from a semantic region of the input document image. An important feature of UDoc is that it learns a generic representation by making use of three self-supervised losses, encouraging the representation to model sentences, learn similarities, and align modalities. Extensive empirical analysis demonstrates that the pretraining procedure learns better joint representations and leads to improvements in downstream tasks.

Finding Discriminative Filters for Specific Degradations in Blind Super-Resolution

- Liangbin Xie, Xintao Wang, Chao Dong, Zhongang Qi, Ying Shan
- abstract@[open-review](#): Recent blind super-resolution (SR) methods typically consist of two branches, one for degradation prediction and the other for conditional restoration. However, our experiments show that a one-branch network can achieve comparable performance to the two-branch scheme. Then we wonder: how can one-branch networks automatically learn to distinguish degradations? To find the answer, we propose a new diagnostic tool -- Filter Attribution method based on Integral Gradient (FAIG). Unlike previous integral gradient methods, our FAIG aims at finding the most discriminative filters instead of input pixels/features for degradation removal in blind SR networks. With the discovered filters, we further develop a simple yet effective method to predict the degradation of an input image. Based on FAIG, we show that, in one-branch blind SR networks, 1) we could find a very small number of (1%) discriminative filters for each specific degradation; 2) The weights, locations and connections of the discovered filters are all important to determine the specific network function. 3) The task of degradation prediction can be implicitly realized by these discriminative filters without explicit supervised learning. Our findings can not only help us better understand network behaviors inside one-branch blind SR networks, but also provide guidance on designing more efficient architectures and diagnosing networks for blind SR.

Counterfactual Explanations Can Be Manipulated

- Dylan Slack, Anna Hilgard, Himabindu Lakkaraju, Sameer Singh
- abstract@[open-review](#): Counterfactual explanations are emerging as an attractive option for providing recourse to individuals adversely impacted by algorithmic decisions. As they are deployed in critical applications (e.g. law enforcement, financial lending), it becomes important to ensure that we clearly understand the vulnerabilities of these methods and find ways to address them. However, there is little understanding of the vulnerabilities and shortcomings of counterfactual explanations. In this work, we introduce the first framework that describes the vulnerabilities of counterfactual explanations and shows how they can be manipulated. More specifically, we show counterfactual explanations may converge to drastically different counterfactuals under a small perturbation indicating they are not robust. Leveraging this insight, we introduce a novel objective to train seemingly fair models where counterfactual explanations find much lower cost recourse under a slight perturbation. We describe how these models can unfairly provide low-cost recourse for specific subgroups in the data while appearing fair to auditors. We perform experiments on loan and violent crime prediction data sets where certain subgroups achieve up to 20x lower cost recourse under the perturbation. These results raise concerns regarding the dependability of current counterfactual explanation techniques, which we hope will inspire investigations in robust counterfactual explanations.

From Canonical Correlation Analysis to Self-supervised Graph Neural Networks

- Hengrui Zhang, Qitian Wu, Junchi Yan, David Wipf, Philip S Yu
- abstract@[open-review](#): We introduce a conceptually simple yet effective model for self-supervised representation learning with graph data. It follows the previous methods that generate two views of an input graph through data augmentation. However, unlike contrastive methods that focus on instance-level discrimination, we optimize an innovative feature-level objective inspired by classical Canonical Correlation Analysis. Compared with other works, our approach requires none of the parameterized mutual information estimator, additional projector, asymmetric structures, and most importantly, negative samples which can be costly. We show that the new objective essentially 1) aims at discarding augmentation-variant information by learning invariant

representations, and 2) can prevent degenerated solutions by decorrelating features in different dimensions. Our theoretical analysis further provides an understanding for the new objective which can be equivalently seen as an instantiation of the Information Bottleneck Principle under the self-supervised setting. Despite its simplicity, our method performs competitively on seven public graph datasets.

[BAST: Bayesian Additive Regression Spanning Trees for Complex Constrained Domain](#)

- Zhao Tang Luo, Huiyan Sang, Bani Mallick
- abstract@[open-review](#): Nonparametric regression on complex domains has been a challenging task as most existing methods, such as ensemble models based on binary decision trees, are not designed to account for intrinsic geometries and domain boundaries. This article proposes a Bayesian additive regression spanning trees (BAST) model for nonparametric regression on manifolds, with an emphasis on complex constrained domains or irregularly shaped spaces embedded in Euclidean spaces. Our model is built upon a random spanning tree manifold partition model as each weak learner, which is capable of capturing any irregularly shaped spatially contiguous partitions while respecting intrinsic geometries and domain boundary constraints. Utilizing many nice properties of spanning tree structures, we design an efficient Bayesian inference algorithm. Equipped with a soft prediction scheme, BAST is demonstrated to significantly outperform other competing methods in simulation experiments and in an application to the chlorophyll data in Aral Sea, due to its strong local adaptivity to different levels of smoothness.

[Hyperbolic Busemann Learning with Ideal Prototypes](#)

- Mina Ghadimi Atigh, Martin Keller-Ressel, Pascal Mettes
- abstract@[open-review](#): Hyperbolic space has become a popular choice of manifold for representation learning of various datatypes from tree-like structures and text to graphs. Building on the success of deep learning with prototypes in Euclidean and hyperspherical spaces, a few recent works have proposed hyperbolic prototypes for classification. Such approaches enable effective learning in low-dimensional output spaces and can exploit hierarchical relations amongst classes, but require privileged information about class labels to position the hyperbolic prototypes. In this work, we propose Hyperbolic Busemann Learning. The main idea behind our approach is to position prototypes on the ideal boundary of the Poincaré ball, which does not require prior label knowledge. To be able to compute proximities to ideal prototypes, we introduce the penalised Busemann loss. We provide theory supporting the use of ideal prototypes and the proposed loss by proving its equivalence to logistic regression in the one-dimensional case. Empirically, we show that our approach provides a natural interpretation of classification confidence, while outperforming recent hyperspherical and hyperbolic prototype approaches.

[Backward-Compatible Prediction Updates: A Probabilistic Approach](#)

- Frederik Träuble, Julius von Kügelgen, Matthäus Kleindessner, Francesco Locatello, Bernhard Schölkopf, Peter Gehler
- abstract@[open-review](#): When machine learning systems meet real world applications, accuracy is only one of several requirements. In this paper, we assay a complementary perspective originating from the increasing availability of pre-trained and regularly improving state-of-the-art models. While new improved models develop at a fast pace, downstream tasks vary more slowly or stay constant. Assume that we have a large unlabelled data set for which we want to maintain accurate predictions. Whenever a new and presumably better ML models becomes available, we encounter two problems: (i) given a limited budget, which data points should be re-evaluated using the new model?; and (ii) if the new predictions differ from the current ones, should we update? Problem (i) is about compute cost, which matters for very large data sets and models. Problem (ii) is about maintaining consistency of the predictions, which can be highly relevant for downstream applications; our demand is to avoid negative flips, i.e., changing correct to incorrect predictions. In this paper, we formalize the Prediction Update Problem and present an efficient probabilistic approach as answer to the above questions. In extensive experiments on standard classification benchmark data sets, we show that our method outperforms alternative strategies along key metrics for backward-compatible prediction updates.

[Truncated Marginal Neural Ratio Estimation](#)

- Benjamin K Miller, Alex Cole, Patrick Forré, Gilles Louppe, Christoph Weniger
- abstract@[open-review](#): Parametric stochastic simulators are ubiquitous in science, often featuring high-dimensional input parameters and/or an intractable likelihood. Performing Bayesian parameter inference in this context can be challenging. We present a neural simulation-based inference algorithm which simultaneously offers simulation efficiency and fast empirical posterior testability, which is unique among modern algorithms. Our approach is simulation efficient by simultaneously estimating low-dimensional marginal posteriors instead of the joint posterior and by proposing simulations targeted to an observation of interest via a prior suitably truncated by an indicator function. Furthermore, by estimating a locally amortized posterior our algorithm enables efficient empirical tests of the robustness of the inference results. Since scientists cannot access the ground truth, these tests are necessary for trusting inference in real-world applications. We perform experiments on a marginalized version of the simulation-based inference benchmark and two complex and narrow posteriors, highlighting the simulator efficiency of our algorithm as well as the quality of the estimated marginal posteriors.

[ReAct: Out-of-distribution Detection With Rectified Activations](#)

- Yiyou Sun, Chuan Guo, Yixuan Li
- abstract@[open-review](#): Out-of-distribution (OOD) detection has received much attention lately due to its practical importance in enhancing the safe deployment of neural networks. One of the primary challenges is that models often produce highly confident predictions on OOD data, which undermines the driving principle in OOD detection that the model should only be confident about in-distribution samples. In this work, we propose ReAct—a simple and effective technique for reducing model overconfidence on OOD data. Our method is motivated by novel analysis on internal activations of neural networks, which displays highly distinctive signature patterns for OOD distributions. Our method can generalize effectively to different network architectures and different OOD detection scores. We empirically demonstrate that ReAct achieves competitive detection performance on a comprehensive suite of benchmark datasets, and give theoretical explication for our method’s efficacy. On the ImageNet benchmark, ReAct reduces the false positive rate (FPR95) by 25.05% compared to the previous best method.

[Non-local Latent Relation Distillation for Self-Adaptive 3D Human Pose Estimation](#)

- Jogendra Nath Kundu, Siddharth Seth, Anirudh Jamkhandi, Pradyumna YM, Varun Jampani, Anirban Chakraborty, Venkatesh Babu R
- abstract@[open-review](#): Available 3D human pose estimation approaches leverage different forms of strong (2D/3D pose) or weak (multi-view or depth) paired supervision. Barring synthetic or in-studio domains, acquiring such supervision for each new target environment is highly inconvenient. To this end, we cast 3D pose learning as a self-supervised adaptation problem that aims to transfer the task knowledge from a labeled source domain to a completely unpaired target. We propose to infer image-to-pose via two explicit mappings viz. image-to-latent and latent-to-pose where the latter is a pre-learned decoder obtained from a prior-enforcing generative adversarial auto-encoder. Next, we introduce relation distillation as a means to align the unpaired cross-modal samples i.e., the unpaired target videos and unpaired 3D pose sequences. To this end, we propose a new set of non-local relations in order to characterize long-range latent pose interactions, unlike general contrastive relations where positive couplings are limited to a local neighborhood structure. Further, we provide an objective way to quantify non-localness in order to select the most effective relation set. We evaluate different self-adaptation settings and demonstrate state-of-the-art 3D human pose estimation performance on standard benchmarks.

[Fast Training of Neural Lumigraph Representations using Meta Learning](#)

- Alexander Bergman, Petr Kellnhofer, Gordon Wetzstein
- abstract@[open-review](#): Novel view synthesis is a long-standing problem in machine learning and computer vision. Significant progress has recently been made in developing neural scene representations and rendering techniques that synthesize photorealistic images from arbitrary views. These representations, however, are extremely slow to train and often also slow to render. Inspired by neural variants of image-based rendering, we develop a new neural rendering approach with the goal of quickly learning a high-quality representation which can also be rendered in real-time. Our approach, MetaNLR++, accomplishes this by using a unique combination of a neural shape representation and 2D CNN-based image feature extraction, aggregation, and re-projection. To push representation convergence times down to minutes, we leverage meta learning to learn neural shape and image feature priors which accelerate training. The optimized shape and image features can then be extracted using traditional graphics techniques and rendered in real time. We show that MetaNLR++ achieves similar or better novel view synthesis results in a fraction of the time that competing methods require.

[Analytical Study of Momentum-Based Acceleration Methods in Paradigmatic High-Dimensional Non-Convex Problems](#)

- Stefano Sarao Mannelli, Pierfrancesco Urbani
- abstract@[open-review](#): The optimization step in many machine learning problems rarely relies on vanilla gradient descent but it is common practice to use momentum-based accelerated methods. Despite these algorithms being widely applied to arbitrary loss functions, their behaviour in generically non-convex, high dimensional landscapes is poorly understood. In this work, we use dynamical mean field theory techniques to describe analytically the average dynamics of these methods in a prototypical non-convex model: the (spiked) matrix-tensor model. We derive a closed set of equations that describe the behaviour of heavy-ball momentum and Nesterov acceleration in the infinite dimensional limit. By numerical integration of these equations, we observe that these methods speed up the dynamics but do not improve the algorithmic threshold with respect to gradient descent in the spiked model.

[Multimodal Few-Shot Learning with Frozen Language Models](#)

- Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, S. M. Ali Eslami, Oriol Vinyals, Felix Hill
- abstract@[open-review](#): When trained at sufficient scale, auto-regressive language models exhibit the notable ability to learn a new language task after being prompted with just a few examples. Here, we present a simple, yet effective, approach for transferring this few-shot learning ability to a multimodal setting (vision and language). Using aligned image and caption data, we train a vision encoder to represent each image as a sequence of continuous embeddings, such that a pre-trained, frozen language model presented with this prefix generates the appropriate caption. The resulting system is a multimodal few-shot learner, with the surprising ability to learn a variety of new tasks when conditioned on examples, represented as a sequence of any number of interleaved image and text embeddings. We demonstrate that it can rapidly learn words for new objects and novel visual categories, do visual question-answering with only a handful of examples, and make use of outside knowledge, by measuring a single model on a variety of established and new benchmarks.

[Approximating the Permanent with Deep Rejection Sampling](#)

- Juha Harviainen, Antti Röyskö, Mikko Koivisto
- abstract@[open-review](#): We present a randomized approximation scheme for the permanent of a matrix with nonnegative entries. Our scheme extends a recursive rejection sampling method of Huber and Law (SODA 2008) by replacing the permanent upper bound with a linear combination of the subproblem bounds at a moderately large depth of the recursion tree. This method, we call deep rejection sampling, is empirically shown to outperform the basic, depth-zero variant, as well as a related method by Kuck et al. (NeurIPS 2019). We analyze the expected running time of the scheme on random $\{0, 1\}$ -matrices where each entry is independently 1 with probability p . Our bound is superior to a previous one for p less than $1/5$, matching another bound that was only known to hold when every row and column has density exactly p .

[Revisiting Model Stitching to Compare Neural Representations](#)

- Yamini Bansal, Preetum Nakkiran, Boaz Barak
- abstract@[open-review](#): We revisit and extend model stitching (Lenc & Vedaldi 2015) as a methodology to study the internal representations of neural networks. Given two trained and frozen models A and B , we consider a "stitched model" formed by connecting the bottom-layers of A to the top-layers of B , with a simple trainable layer between them. We argue that model stitching is a powerful and perhaps under-appreciated tool, which reveals aspects of representations that measures such as centered kernel alignment (CKA) cannot. Through extensive experiments, we use model stitching to obtain quantitative verifications for intuitive statements such as "good networks learn similar representations", by demonstrating that good networks of the same architecture, but trained in very different ways (eg: supervised vs. self-supervised learning), can be stitched to each other without drop in performance. We also give evidence for the intuition that "more is better" by showing that representations learnt with (1) more data, (2) bigger width, or (3) more training time can be "plugged in" to weaker models to improve performance. Finally, our experiments reveal a new structural property of SGD which we call "stitching connectivity", akin to mode-connectivity: typical minima reached by SGD are all "stitching-connected" to each other.

[AugMax: Adversarial Composition of Random Augmentations for Robust Training](#)

- Haotao Wang, Chaowei Xiao, Jean Kossaifi, Zhiding Yu, Anima Anandkumar, Zhangyang Wang
- abstract@[open-review](#): Data augmentation is a simple yet effective way to improve the robustness of deep neural networks (DNNs). Diversity and hardness are two complementary dimensions of data augmentation to achieve robustness. For example, AugMix explores random compositions of a diverse set of augmentations to enhance broader coverage, while adversarial training generates adversarially hard samples to spot the weakness. Motivated by this, we propose a data augmentation framework, termed AugMax, to unify the two aspects of diversity and hardness. AugMax first randomly samples multiple augmentation operators and then learns an adversarial mixture of the selected operators. Being a stronger form of data augmentation, AugMax leads to a significantly augmented input distribution which makes model training more challenging. To solve this problem, we further design a disentangled normalization module, termed DuBIN (Dual-Batch-and-Instance Normalization), that disentangles the instance-wise feature heterogeneity arising from AugMax. Experiments show that AugMax-DuBIN leads to significantly improved out-of-distribution robustness, outperforming prior arts by 3.03%, 3.49%, 1.82% and 0.71% on CIFAR10-C, CIFAR100-C, Tiny ImageNet-C and ImageNet-C. Codes and pretrained models are available: <https://github.com/VITA-Group/AugMax>.

[Habitat 2.0: Training Home Assistants to Rearrange their Habitat](#)

- Andrew Szot, Alexander Clegg, Eric Undersander, Erik Wijmans, Yili Zhao, John Turner, Noah Maestre, Mustafa Mukadam, Devendra Singh Chaplot, Oleksandr MakSYMets, Aaron Gokaslan, Vladimír Vondruš, Sameer Dharur, Franziska Meier, Wojciech Galuba, Angel Chang, Zsolt Kira, Vladlen Koltun, Jitendra Malik, Manolis Savva, Dhruv Batra
- abstract@[open-review](#): We introduce Habitat 2.0 (H2.0), a simulation platform for training virtual robots in interactive 3D environments and complex physics-enabled scenarios. We make comprehensive contributions to all levels of the embodied AI stack – data, simulation, and benchmark tasks. Specifically, we present: (i) ReplicaCAD: an artist-authored, annotated, reconfigurable 3D dataset of apartments (matching real spaces) with articulated objects (e.g. cabinets and drawers that can open/close); (ii) H2.0: a high-performance physics-enabled 3D simulator with speeds exceeding 25,000 simulation steps per second (850x real-time) on an 8-GPU node, representing 100x speed-ups over prior work; and, (iii) Home Assistant Benchmark (HAB): a suite of common tasks for assistive robots (tidy the house, stock groceries, set the table) that test a range of mobile manipulation capabilities. These large-scale engineering contributions allow us to systematically compare deep reinforcement learning (RL) at scale and classical sense-plan-act (SPA) pipelines in long-horizon structured tasks, with an emphasis on generalization to new objects, receptacles, and layouts. We find that (1) flat RL

policies struggle on HAB compared to hierarchical ones; (2) a hierarchy with independent skills suffers from ‘hand-off problems’, and (3) SPA pipelines are more brittle than RL policies.

[Time Discretization-Invariant Safe Action Repetition for Policy Gradient Methods](#)

- Seohong Park, Jaekyeom Kim, Gunhee Kim
- abstract@[open-review](#): In reinforcement learning, continuous time is often discretized by a time scale δ , to which the resulting performance is known to be highly sensitive. In this work, we seek to find a δ -invariant algorithm for policy gradient (PG) methods, which performs well regardless of the value of δ . We first identify the underlying reasons that cause PG methods to fail as $\delta \rightarrow 0$, proving that the variance of the PG estimator can diverge to infinity in stochastic environments under a certain assumption of stochasticity. While durative actions or action repetition can be employed to have δ -invariance, previous action repetition methods cannot immediately react to unexpected situations in stochastic environments. We thus propose a novel δ -invariant method named Safe Action Repetition (SAR) applicable to any existing PG algorithm. SAR can handle the stochasticity of environments by adaptively reacting to changes in states during action repetition. We empirically show that our method is not only δ -invariant but also robust to stochasticity, outperforming previous δ -invariant approaches on eight MuJoCo environments with both deterministic and stochastic settings. Our code is available at <https://vision.snu.ac.kr/projects/sar>.

[Meta-Learning Reliable Priors in the Function Space](#)

- Jonas Rothfuss, Dominique Heyn, jinfan Chen, Andreas Krause
- abstract@[open-review](#): Meta-Learning promises to enable more data-efficient inference by harnessing previous experience from related learning tasks. While existing meta-learning methods help us to improve the accuracy of our predictions in face of data scarcity, they fail to supply reliable uncertainty estimates, often being grossly overconfident in their predictions. Addressing these shortcomings, we introduce a novel meta-learning framework, called F-PACOH, that treats meta-learned priors as stochastic processes and performs meta-level regularization directly in the function space. This allows us to directly steer the probabilistic predictions of the meta-learner towards high epistemic uncertainty in regions of insufficient meta-training data and, thus, obtain well-calibrated uncertainty estimates. Finally, we showcase how our approach can be integrated with sequential decision making, where reliable uncertainty quantification is imperative. In our benchmark study on meta-learning for Bayesian Optimization (BO), F-PACOH significantly outperforms all other meta-learners and standard baselines. Even in a challenging lifelong BO setting, where optimization tasks arrive one at a time and the meta-learner needs to build up informative prior knowledge incrementally, our proposed method demonstrates strong positive transfer.

[VoiceMixer: Adversarial Voice Style Mixup](#)

- Sang-Hoon Lee, Ji-Hoon Kim, Hyunseung Chung, Seong-Whan Lee
- abstract@[open-review](#): Although recent advances in voice conversion have shown significant improvement, there still remains a gap between the converted voice and target voice. A key factor that maintains this gap is the insufficient decomposition of content and voice style from the source speech. This insufficiency leads to the converted speech containing source speech style or losing source speech content. In this paper, we present VoiceMixer which can effectively decompose and transfer voice style through a novel information bottleneck and adversarial feedback. With self-supervised representation learning, the proposed information bottleneck can decompose the content and style with only a small loss of content information. Also, for adversarial feedback of each information, the discriminator is decomposed into content and style discriminator with self-supervision, which enable our model to achieve better generalization to the voice style of the converted speech. The experimental results show the superiority of our model in disentanglement and transfer performance, and improve audio quality by preserving content information.

[Predicting What You Already Know Helps: Provable Self-Supervised Learning](#)

- Jason D. Lee, Qi Lei, Nikunj Saunshi, JIACHENG ZHUO
- abstract@[open-review](#): Self-supervised representation learning solves auxiliary prediction tasks (known as pretext tasks), that do not require labeled data, to learn semantic representations. These pretext tasks are created solely using the input features, such as predicting a missing image patch, recovering the color channels of an image from context, or predicting missing words, yet predicting this \textit{known} information helps in learning representations effective for downstream prediction tasks. This paper posits a mechanism based on approximate conditional independence to formalize how solving certain pretext tasks can learn representations that provably decrease the sample complexity of downstream supervised tasks. Formally, we quantify how the approximate independence between the components of the pretext task (conditional on the label and latent variables) allows us to learn representations that can solve the downstream task with drastically reduced sample complexity by just training a linear layer on top of the learned representation.

[Oracle Complexity in Nonsmooth Nonconvex Optimization](#)

- Guy Kornowski, Ohad Shamir
- abstract@[open-review](#): It is well-known that given a smooth, bounded-from-below, and possibly nonconvex function, standard gradient-based methods can find ϵ -stationary points (with gradient norm less than ϵ) in $\mathcal{O}(1/\epsilon^2)$ iterations. However, many important nonconvex optimization problems, such as those associated with training modern neural networks, are inherently not smooth, making these results inapplicable. In this paper, we study nonsmooth nonconvex optimization from an oracle complexity viewpoint, where the algorithm is assumed to be given access only to local information about the function at various points. We provide two main results (under mild assumptions): First, we consider the problem of getting \textit{near} ϵ -stationary points. This is perhaps the most natural relaxation of \textit{finding} ϵ -stationary points, which is impossible in the nonsmooth nonconvex case. We prove that this relaxed goal cannot be achieved efficiently, for any distance and ϵ smaller than some constants. Our second result deals with the possibility of tackling nonsmooth nonconvex optimization by reduction to smooth optimization: Namely, applying smooth optimization methods on a smooth approximation of the objective function. For this approach, we prove an inherent trade-off between oracle complexity and smoothness: On the one hand, smoothing a nonsmooth nonconvex function can be done very efficiently (e.g., by randomized smoothing), but with dimension-dependent factors in the smoothness parameter, which can strongly affect iteration complexity when plugging into standard smooth optimization methods. On the other hand, these dimension factors can be eliminated with suitable smoothing methods, but only by making the oracle complexity of the smoothing process exponentially large.

[CentripetalText: An Efficient Text Instance Representation for Scene Text Detection](#)

- Tao Sheng, Jie Chen, Zhouhui Lian
- abstract@[open-review](#): Scene text detection remains a grand challenge due to the variation in text curvatures, orientations, and aspect ratios. One of the hardest problems in this task is how to represent text instances of arbitrary shapes. Although many methods have been proposed to model irregular texts in a flexible manner, most of them lose simplicity and robustness. Their complicated post-processings and the regression under Dirac delta distribution undermine the detection performance and the generalization ability. In this paper, we propose an efficient text instance representation named CentripetalText (CT), which decomposes text instances into the combination of text kernels and centripetal shifts. Specifically, we utilize the centripetal shifts to implement pixel aggregation, guiding the external text pixels to the internal text kernels. The relaxation operation is integrated into the dense regression for centripetal shifts, allowing the correct prediction in a range instead of a specific value. The convenient reconstruction of text contours and the tolerance of prediction errors in our method guarantee the high detection accuracy and the fast inference speed, respectively. Besides, we shrink our text detector into a proposal generation module, namely CentripetalText Proposal Network (CPN), replacing Segmentation Proposal Network (SPN) in Mask TextSpotter v3 and producing more accurate proposals. To validate the effectiveness of our method, we conduct experiments on several commonly

used scene text benchmarks, including both curved and multi-oriented text datasets. For the task of scene text detection, our approach achieves superior or competitive performance compared to other existing methods, e.g., F-measure of 86.3% at 40.0 FPS on Total-Text, F-measure of 86.1% at 34.8 FPS on MSRA-TD500, etc. For the task of end-to-end scene text recognition, our method outperforms Mask TextSpotter v3 by 1.1% in F-measure on Total-Text.

[Learning to Select Exogenous Events for Marked Temporal Point Process](#)

- Ping Zhang, Rishabh Iyer, Ashish Tendulkar, Gaurav Aggarwal, Abir De
- abstract@[open-review](#): Marked temporal point processes (MTPPs) have emerged as a powerful modeling tool for a wide variety of applications which are characterized using discrete events localized in continuous time. In this context, the events are of two types: endogenous events which occur due to the influence of the previous events and exogenous events which occur due to the effect of the externalities. However, in practice, the events do not come with endogenous or exogenous labels. To this end, our goal in this paper is to identify the set of exogenous events from a set of unlabelled events. To do so, we first formulate the parameter estimation problem in conjunction with exogenous event set selection problem and show that this problem is NP hard. Next, we prove that the underlying objective is a monotone and α -submodular set function, with respect to the candidate set of exogenous events. Such a characterization subsequently allows us to use a stochastic greedy algorithm which was originally proposed in [\cite{greedy}](#) for submodular maximization. However, we show that it also admits an approximation guarantee for maximizing α -submodular set function, even when the learning algorithm provides an imperfect estimate of the trained parameters. Finally, our experiments with synthetic and real data show that our method performs better than the existing approaches built upon superposition of endogenous and exogenous MTPPs.

[DRIVE: One-bit Distributed Mean Estimation](#)

- Shay Vargaftik, Ran Ben-Basat, Amit Portnoy, Gal Mendelson, Yaniv Ben-Itzhak, Michael Mitzenmacher
- abstract@[open-review](#): We consider the problem where n clients transmit d -dimensional real-valued vectors using $d(1+o(1))$ bits each, in a manner that allows the receiver to approximately reconstruct their mean. Such compression problems naturally arise in distributed and federated learning. We provide novel mathematical results and derive computationally efficient algorithms that are more accurate than previous compression techniques. We evaluate our methods on a collection of distributed and federated learning tasks, using a variety of datasets, and show a consistent improvement over the state of the art.

[Learning Space Partitions for Path Planning](#)

- Kevin Yang, Tianjun Zhang, Chris Cummins, Brandon Cui, Benoit Steiner, Linnan Wang, Joseph E. Gonzalez, Dan Klein, Yuandong Tian
- abstract@[open-review](#): Path planning, the problem of efficiently discovering high-reward trajectories, often requires optimizing a high-dimensional and multimodal reward function. Popular approaches like CEM and CMA-ES greedily focus on promising regions of the search space and may get trapped in local maxima. DOO and VOOT balance exploration and exploitation, but use space partitioning strategies independent of the reward function to be optimized. Recently, LaMCTS empirically learns to partition the search space in a reward-sensitive manner for black-box optimization. In this paper, we develop a novel formal regret analysis for when and why such an adaptive region partitioning scheme works. We also propose a new path planning method LaP3 which improves the function value estimation within each sub-region, and uses a latent representation of the search space. Empirically, LaP3 outperforms existing path planning methods in 2D navigation tasks, especially in the presence of difficult-to-escape local optima, and shows benefits when plugged into the planning components of model-based RL such as PETS. These gains transfer to highly multimodal real-world tasks, where we outperform strong baselines in compiler phase ordering by up to 39% on average across 9 tasks, and in molecular design by up to 0.4 on properties on a 0-1 scale. Code is available at <https://github.com/yangkevin2/neurips2021-lap3>.

[Progressive Feature Interaction Search for Deep Sparse Network](#)

- Chen Gao, Yinfeng Li, Quanming Yao, Depeng Jin, Yong Li
- abstract@[open-review](#): Deep sparse networks (DSNs), of which the crux is exploring the high-order feature interactions, have become the state-of-the-art on the prediction task with high-sparsity features. However, these models suffer from low computation efficiency, including large model size and slow model inference, which largely limits these models' application value. In this work, we approach this problem with neural architecture search by automatically searching the critical component in DSNs, the feature-interaction layer. We propose a distilled search space to cover the desired architectures with fewer parameters. We then develop a progressive search algorithm for efficient search on the space and well capture the order-priority property in sparse prediction tasks. Experiments on three real-world benchmark datasets show promising results of PROFIT in both accuracy and efficiency. Further studies validate the feasibility of our designed search space and search algorithm.

[Local Explanation of Dialogue Response Generation](#)

- Yi-Lin Tuan, Connor Pryor, Wenhui Chen, Lise Getoor, William Yang Wang
- abstract@[open-review](#): In comparison to the interpretation of classification models, the explanation of sequence generation models is also an important problem, however it has seen little attention. In this work, we study model-agnostic explanations of a representative text generation task -- dialogue response generation. Dialog response generation is challenging with its open-ended sentences and multiple acceptable responses. To gain insights into the reasoning process of a generation model, we propose a new method, local explanation of response generation (LERG) that regards the explanations as the mutual interaction of segments in input and output sentences. LERG views the sequence prediction as uncertainty estimation of a human response and then creates explanations by perturbing the input and calculating the certainty change over the human response. We show that LERG adheres to desired properties of explanations for text generation including unbiased approximation, consistency and cause identification. Empirically, our results show that our method consistently improves other widely used methods on proposed automatic- and human- evaluation metrics for this new task by \$4.4\%\$-\$12.8\%\$. Our analysis demonstrates that LERG can extract both explicit and implicit relations between input and output segments.

[Scalable Inference in SDEs by Direct Matching of the Fokker–Planck–Kolmogorov Equation](#)

- Arno Solin, Ella Tamir, Prakhar Verma
- abstract@[open-review](#): Simulation-based techniques such as variants of stochastic Runge–Kutta are the de facto approach for inference with stochastic differential equations (SDEs) in machine learning. These methods are general-purpose and used with parametric and non-parametric models, and neural SDEs. Stochastic Runge–Kutta relies on the use of sampling schemes that can be inefficient in high dimensions. We address this issue by revisiting the classical SDE literature and derive direct approximations to the (typically intractable) Fokker–Planck–Kolmogorov equation by matching moments. We show how this workflow is fast, scales to high-dimensional latent spaces, and is applicable to scarce-data applications, where a non-parametric SDE with a driving Gaussian process velocity field specifies the model.

[The Complexity of Bayesian Network Learning: Revisiting the Superstructure](#)

- Robert Ganian, Viktoriia Korchemna
- abstract@[open-review](#): We investigate the parameterized complexity of Bayesian Network Structure Learning (BNSL), a classical problem that has received significant attention in empirical but also purely theoretical studies. We follow up on previous works that have analyzed the complexity of BNSL w.r.t. the so-called superstructure of the input. While known results imply that BNSL is unlikely to be fixed-parameter tractable even when parameterized

by the size of a vertex cover in the superstructure, here we show that a different kind of parameterization - notably by the size of a feedback edge set - yields fixed-parameter tractability. We proceed by showing that this result can be strengthened to a localized version of the feedback edge set, and provide corresponding lower bounds that complement previous results to provide a complexity classification of BNSL w.r.t. virtually all well-studied graph parameters. We then analyze how the complexity of BNSL depends on the representation of the input. In particular, while the bulk of past theoretical work on the topic assumed the use of the so-called non-zero representation, here we prove that if an additive representation can be used instead then BNSL becomes fixed-parameter tractable even under significantly milder restrictions to the superstructure, notably when parameterized by the treewidth alone. Last but not least, we show how our results can be extended to the closely related problem of Polytree Learning.

[Fast Tucker Rank Reduction for Non-Negative Tensors Using Mean-Field Approximation](#)

- Kazu Ghalamkari, Mahito Sugiyama
- abstract@[open-review](#): We present an efficient low-rank approximation algorithm for non-negative tensors. The algorithm is derived from our two findings: First, we show that rank-1 approximation for tensors can be viewed as a mean-field approximation by treating each tensor as a probability distribution. Second, we theoretically provide a sufficient condition for distribution parameters to reduce Tucker ranks of tensors; interestingly, this sufficient condition can be achieved by iterative application of the mean-field approximation. Since the mean-field approximation is always given as a closed formula, our findings lead to a fast low-rank approximation algorithm without using a gradient method. We empirically demonstrate that our algorithm is faster than the existing non-negative Tucker rank reduction methods and achieves competitive or better approximation of given tensors.

[Learning Stochastic Majority Votes by Minimizing a PAC-Bayes Generalization Bound](#)

- Valentina Zantedeschi, Paul Viallard, Emilie Morvant, Rémi Emonet, Amaury Habrard, Pascal Germain, Benjamin Guedj
- abstract@[open-review](#): We investigate a stochastic counterpart of majority votes over finite ensembles of classifiers, and study its generalization properties. While our approach holds for arbitrary distributions, we instantiate it with Dirichlet distributions: this allows for a closed-form and differentiable expression for the expected risk, which then turns the generalization bound into a tractable training objective. The resulting stochastic majority vote learning algorithm achieves state-of-the-art accuracy and benefits from (non-vacuous) tight generalization bounds, in a series of numerical experiments when compared to competing algorithms which also minimize PAC-Bayes objectives -- both with uninformed (data-independent) and informed (data-dependent) priors.

[Numerical influence of ReLU'\(0\) on backpropagation](#)

- David Bertoin, Jérôme Bolte, Sébastien Gerchinovitz, Edouard Pauwels
- abstract@[open-review](#): In theory, the choice of ReLU(0) in [0, 1] for a neural network has a negligible influence both on backpropagation and training. Yet, in the real world, 32 bits default precision combined with the size of deep learning problems makes it a hyperparameter of training methods. We investigate the importance of the value of ReLU'(0) for several precision levels (16, 32, 64 bits), on various networks (fully connected, VGG, ResNet) and datasets (MNIST, CIFAR10, SVHN, ImageNet). We observe considerable variations of backpropagation outputs which occur around half of the time in 32 bits precision. The effect disappears with double precision, while it is systematic at 16 bits. For vanilla SGD training, the choice ReLU'(0) = 0 seems to be the most efficient. For our experiments on ImageNet the gain in test accuracy over ReLU'(0) = 1 was more than 10 points (two runs). We also evidence that reconditioning approaches as batch-norm or ADAM tend to buffer the influence of ReLU'(0)'s value. Overall, the message we convey is that algorithmic differentiation of nonsmooth problems potentially hides parameters that could be tuned advantageously.

[A Contrastive Learning Approach for Training Variational Autoencoder Priors](#)

- Jyoti Aneja, Alex Schwing, Jan Kautz, Arash Vahdat
- abstract@[open-review](#): Variational autoencoders (VAEs) are one of the powerful likelihood-based generative models with applications in many domains. However, they struggle to generate high-quality images, especially when samples are obtained from the prior without any tempering. One explanation for VAEs' poor generative quality is the prior hole problem: the prior distribution fails to match the aggregate approximate posterior. Due to this mismatch, there exist areas in the latent space with high density under the prior that do not correspond to any encoded image. Samples from those areas are decoded to corrupted images. To tackle this issue, we propose an energy-based prior defined by the product of a base prior distribution and a reweighting factor, designed to bring the base closer to the aggregate posterior. We train the reweighting factor by noise contrastive estimation, and we generalize it to hierarchical VAEs with many latent variable groups. Our experiments confirm that the proposed noise contrastive priors improve the generative performance of state-of-the-art VAEs by a large margin on the MNIST, CIFAR-10, CelebA 64, and CelebA HQ 256 datasets. Our method is simple and can be applied to a wide variety of VAEs to improve the expressivity of their prior distribution.

[What training reveals about neural network complexity](#)

- Andreas Loukas, Marinos Poiitis, Stefanie Jegelka
- abstract@[open-review](#): This work explores the Benevolent Training Hypothesis (BTH) which argues that the complexity of the function a deep neural network (NN) is learning can be deduced by its training dynamics. Our analysis provides evidence for BTH by relating the NN's Lipschitz constant at different regions of the input space with the behavior of the stochastic training procedure. We first observe that the Lipschitz constant close to the training data affects various aspects of the parameter trajectory, with more complex networks having a longer trajectory, bigger variance, and often veering further from their initialization. We then show that NNs whose 1st layer bias is trained more steadily (i.e., slowly and with little variation) have bounded complexity even in regions of the input space that are far from any training point. Finally, we find that steady training with Dropout implies a training- and data-dependent generalization bound that grows poly-logarithmically with the number of parameters. Overall, our results support the intuition that good training behavior can be a useful bias towards good generalization.

[Class-agnostic Reconstruction of Dynamic Objects from Videos](#)

- Zhongzheng Ren, Xiaoming Zhao, Alex Schwing
- abstract@[open-review](#): We introduce REDO, a class-agnostic framework to REconstruct the Dynamic Objects from RGBD or calibrated videos. Compared to prior work, our problem setting is more realistic yet more challenging for three reasons: 1) due to occlusion or camera settings an object of interest may never be entirely visible, but we aim to reconstruct the complete shape; 2) we aim to handle different object dynamics including rigid motion, non-rigid motion, and articulation; 3) we aim to reconstruct different categories of objects with one unified framework. To address these challenges, we develop two novel modules. First, we introduce a canonical 4D implicit function which is pixel-aligned with aggregated temporal visual cues. Second, we develop a 4D transformation module which captures object dynamics to support temporal propagation and aggregation. We study the efficacy of REDO in extensive experiments on synthetic RGBD video datasets SAIL-VOS 3D and DeformingThings4D++, and on real-world video data 3DPW. We find REDO outperforms state-of-the-art dynamic reconstruction methods by a margin. In ablation studies we validate each developed component.

[Unique sparse decomposition of low rank matrices](#)

- Dian Jin, Xin Bing, Yuqian Zhang

- abstract@[open-review](#): The problem of finding the unique low dimensional decomposition of a given matrix has been a fundamental and recurrent problem in many areas. In this paper, we study the problem of seeking a unique decomposition of a low-rank matrix $Y \in \mathbb{R}^{p \times n}$ that admits a sparse representation. Specifically, we consider $Y = AX \in \mathbb{R}^{p \times n}$ where the matrix $A \in \mathbb{R}^{p \times r}$ has full column rank, with $r < \min\{n, p\}$, and the matrix $X \in \mathbb{R}^{r \times n}$ is element-wise sparse. We prove that this sparse decomposition of Y can be uniquely identified by recovering ground-truth A column by column, up to some intrinsic signed permutation. Our approach relies on solving a nonconvex optimization problem constrained over the unit sphere. Our geometric analysis for the nonconvex optimization landscape shows that any local solution is close to the ground truth solution, and can be recovered by a simple data-driven initialization followed with any second-order descent algorithm. At last, we corroborate these theoretical results with numerical experiments

[Neighborhood Reconstructing Autoencoders](#)

- Yonghyeon LEE, Hyeokjun Kwon, Frank Park
- abstract@[open-review](#): Vanilla autoencoders often produce manifolds that overfit to noisy training data, or have the wrong local connectivity and geometry. Autoencoder regularization techniques, e.g., the denoising autoencoder, have had some success in reducing overfitting, whereas recent graph-based methods that exploit local connectivity information provided by neighborhood graphs have had some success in mitigating local connectivity errors. Neither of these two approaches satisfactorily reduce both overfitting and connectivity errors; moreover, graph-based methods typically involve considerable preprocessing and tuning. To simultaneously address the two issues of overfitting and local connectivity, we propose a new graph-based autoencoder, the Neighborhood Reconstructing Autoencoder (NRAE). Unlike existing graph-based methods that attempt to encode the training data to some prescribed latent space distribution -- one consequence being that only the encoder is the object of the regularization -- NRAE merges local connectivity information contained in the neighborhood graphs with local quadratic approximations of the decoder function to formulate a new neighborhood reconstruction loss. Compared to existing graph-based methods, our new loss function is simple and easy to implement, and the resulting algorithm is scalable and computationally efficient; the only required preprocessing step is the construction of the neighborhood graph. Extensive experiments with standard datasets demonstrate that, compared to existing methods, NRAE improves both overfitting and local connectivity in the learned manifold, in some cases by significant margins. Code for NRAE is available at <https://github.com/Gabe-YHLee/NRAE-public>.

[TopicNet: Semantic Graph-Guided Topic Discovery](#)

- Zhibin Duan, Yi Shi Xu, Bo Chen, Dongsheng Wang, Chaojie Wang, Mingyuan Zhou
- abstract@[open-review](#): Existing deep hierarchical topic models are able to extract semantically meaningful topics from a text corpus in an unsupervised manner and automatically organize them into a topic hierarchy. However, it is unclear how to incorporate prior belief such as knowledge graph to guide the learning of the topic hierarchy. To address this issue, we introduce TopicNet as a deep hierarchical topic model that can inject prior structural knowledge as inductive bias to influence the learning. TopicNet represents each topic as a Gaussian-distributed embedding vector, projects the topics of all layers into a shared embedding space, and explores both the symmetric and asymmetric similarities between Gaussian embedding vectors to incorporate prior semantic hierarchies. With a variational auto-encoding inference network, the model parameters are optimized by minimizing the evidence lower bound and supervised loss via stochastic gradient descent. Experiments on widely used benchmark show that TopicNet outperforms related deep topic models on discovering deeper interpretable topics and mining better document representations.

[\(Almost\) Free Incentivized Exploration from Decentralized Learning Agents](#)

- Chengshuai Shi, Haifeng Xu, Wei Xiong, Cong Shen
- abstract@[open-review](#): Incentivized exploration in multi-armed bandits (MAB) has witnessed increasing interests and many progresses in recent years, where a principal offers bonuses to agents to do explorations on her behalf. However, almost all existing studies are confined to temporary myopic agents. In this work, we break this barrier and study incentivized exploration with multiple and long-term strategic agents, who have more complicated behaviors that often appear in real-world applications. An important observation of this work is that strategic agents' intrinsic needs of learning benefit (instead of harming) the principal's explorations by providing "free pulls". Moreover, it turns out that increasing the population of agents significantly lowers the principal's burden of incentivizing. The key and somewhat surprising insight revealed from our results is that when there are sufficiently many learning agents involved, the exploration process of the principal can be (almost) free. Our main results are built upon three novel components which may be of independent interest: (1) a simple yet provably effective incentive-provision strategy; (2) a carefully crafted best arm identification algorithm for rewards aggregated under unequal confidences; (3) a high-probability finite-time lower bound of UCB algorithms. Experimental results are provided to complement the theoretical analysis.

[Combining Recurrent, Convolutional, and Continuous-time Models with Linear State Space Layers](#)

- Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, Christopher Ré
- abstract@[open-review](#): Recurrent neural networks (RNNs), temporal convolutions, and neural differential equations (NDEs) are popular families of deep learning models for time-series data, each with unique strengths and tradeoffs in modeling power and computational efficiency. We introduce a simple sequence model inspired by control systems that generalizes these approaches while addressing their shortcomings. The Linear State-Space Layer (LSSL) maps a sequence $u \mapsto y$ by simply simulating a linear continuous-time state-space representation $\dot{x} = Ax + Bu, y = Cx + Du$. Theoretically, we show that LSSL models are closely related to the three aforementioned families of models and inherit their strengths. For example, they generalize convolutions to continuous-time, explain common RNN heuristics, and share features of NDEs such as time-scale adaptation. We then incorporate and generalize recent theory on continuous-time memorization to introduce a trainable subset of structured matrices A that endow LSSLs with long-range memory. Empirically, stacking LSSL layers into a simple deep neural network obtains state-of-the-art results across time series benchmarks for long dependencies in sequential image classification, real-world healthcare regression tasks, and speech. On a difficult speech classification task with length-16000 sequences, LSSL outperforms prior approaches by 24 accuracy points, and even outperforms baselines that use hand-crafted features on 100x shorter sequences.

[Revisiting Hilbert-Schmidt Information Bottleneck for Adversarial Robustness](#)

- Zifeng Wang, Tong Jian, Aria Masoomi, Stratis Ioannidis, Jennifer Dy
- abstract@[open-review](#): We investigate the HSIC (Hilbert-Schmidt independence criterion) bottleneck as a regularizer for learning an adversarially robust deep neural network classifier. In addition to the usual cross-entropy loss, we add regularization terms for every intermediate layer to ensure that the latent representations retain useful information for output prediction while reducing redundant information. We show that the HSIC bottleneck enhances robustness to adversarial attacks both theoretically and experimentally. In particular, we prove that the HSIC bottleneck regularizer reduces the sensitivity of the classifier to adversarial examples. Our experiments on multiple benchmark datasets and architectures demonstrate that incorporating an HSIC bottleneck regularizer attains competitive natural accuracy and improves adversarial robustness, both with and without adversarial examples during training. Our code and adversarially robust models are publicly available.

[T-LoHo: A Bayesian Regularization Model for Structured Sparsity and Smoothness on Graphs](#)

- Changwoo Lee, Zhao Tang Luo, Huiyan Sang
- abstract@[open-review](#): Graphs have been commonly used to represent complex data structures. In models dealing with graph-structured data, multivariate parameters may not only exhibit sparse patterns but have structured sparsity and smoothness in the sense that both zero and non-zero parameters tend to

cluster together. We propose a new prior for high-dimensional parameters with graphical relations, referred to as the Tree-based Low-rank Horseshoe (T-LoHo) model, that generalizes the popular univariate Bayesian horseshoe shrinkage prior to the multivariate setting to detect structured sparsity and smoothness simultaneously. The T-LoHo prior can be embedded in many high-dimensional hierarchical models. To illustrate its utility, we apply it to regularize a Bayesian high-dimensional regression problem where the regression coefficients are linked by a graph, so that the resulting clusters have flexible shapes and satisfy the cluster contiguity constraint with respect to the graph. We design an efficient Markov chain Monte Carlo algorithm that delivers full Bayesian inference with uncertainty measures for model parameters such as the number of clusters. We offer theoretical investigations of the clustering effects and posterior concentration results. Finally, we illustrate the performance of the model with simulation studies and a real data application for anomaly detection on a road network. The results indicate substantial improvements over other competing methods such as the sparse fused lasso.

The Utility of Explainable AI in Ad Hoc Human-Machine Teaming

- Rohan Paleja, Muyleng Ghuy, Nadun Ranawaka Arachchige, Reed Jensen, Matthew Gombolay
- abstract@[open-review](#): Recent advances in machine learning have led to growing interest in Explainable AI (xAI) to enable humans to gain insight into the decision-making of machine learning models. Despite this recent interest, the utility of xAI techniques has not yet been characterized in human-machine teaming. Importantly, xAI offers the promise of enhancing team situational awareness (SA) and shared mental model development, which are the key characteristics of effective human-machine teams. Rapidly developing such mental models is especially critical in ad hoc human-machine teaming, where agents do not have a priori knowledge of others' decision-making strategies. In this paper, we present two novel human-subject experiments quantifying the benefits of deploying xAI techniques within a human-machine teaming scenario. First, we show that xAI techniques can support SA ($p < 0.05$). Second, we examine how different SA levels induced via a collaborative AI policy abstraction affect ad hoc human-machine teaming performance. Importantly, we find that the benefits of xAI are not universal, as there is a strong dependence on the composition of the human-machine team. Novices benefit from xAI providing increased SA ($p < 0.05$) but are susceptible to cognitive overhead ($p < 0.05$). On the other hand, expert performance degrades with the addition of xAI-based support ($p < 0.05$), indicating that the cost of paying attention to the xAI outweighs the benefits obtained from being provided additional information to enhance SA. Our results demonstrate that researchers must deliberately design and deploy the right xAI techniques in the right scenario by carefully considering human-machine team composition and how the xAI method augments SA.

Subgoal Search For Complex Reasoning Tasks

- Konrad Czechowski, Tomasz Odrzygóźdż, Marek Zbysiński, Michał Zawalski, Krzysztof Olejnik, Yuhuai Wu, Łukasz Kuciński, Piotr Miłoś
- abstract@[open-review](#): Humans excel in solving complex reasoning tasks through a mental process of moving from one idea to a related one. Inspired by this, we propose Subgoal Search (kSubS) method. Its key component is a learned subgoal generator that produces a diversity of subgoals that are both achievable and closer to the solution. Using subgoals reduces the search space and induces a high-level search graph suitable for efficient planning. In this paper, we implement kSubS using a transformer-based subgoal module coupled with the classical best-first search framework. We show that a simple approach of generating k -th step ahead subgoals is surprisingly efficient on three challenging domains: two popular puzzle games, Sokoban and the Rubik's Cube, and an inequality proving benchmark INT. kSubS achieves strong results including state-of-the-art on INT within a modest computational budget.

MCMC Variational Inference via Uncorrected Hamiltonian Annealing

- Tomas Geffner, Justin Domke
- abstract@[open-review](#): Given an unnormalized target distribution we want to obtain approximate samples from it and a tight lower bound on its (log) normalization constant $\log Z$. Annealed Importance Sampling (AIS) with Hamiltonian MCMC is a powerful method that can be used to do this. Its main drawback is that it uses non-differentiable transition kernels, which makes tuning its many parameters hard. We propose a framework to use an AIS-like procedure with Uncorrected Hamiltonian MCMC, called Uncorrected Hamiltonian Annealing. Our method leads to tight and differentiable lower bounds on $\log Z$. We show empirically that our method yields better performances than other competing approaches, and that the ability to tune its parameters using reparameterization gradients may lead to large performance improvements.

Landmark-RxR: Solving Vision-and-Language Navigation with Fine-Grained Alignment Supervision

- Keji He, Yan Huang, Qi Wu, Jianhua Yang, Dong An, Shuanglin Sima, Liang Wang
- abstract@[open-review](#): In Vision-and-Language Navigation (VLN) task, an agent is asked to navigate inside 3D indoor environments following given instructions. Cross-modal alignment is one of the most critical challenges in VLN because the predicted trajectory needs to match the given instruction accurately. In this paper, we address the cross-modal alignment challenge from the perspective of fine-grain. Firstly, to alleviate weak cross-modal alignment supervision from coarse-grained data, we introduce a human-annotated fine-grained VLN dataset, namely Landmark-RxR. Secondly, to further enhance local cross-modal alignment under fine-grained supervision, we investigate the focal-oriented rewards with soft and hard forms, by focusing on the critical points sampled from fine-grained Landmark-RxR. Moreover, to fully evaluate the navigation process, we also propose a re-initialization mechanism that makes metrics insensitive to difficult points, which can cause the agent to deviate from the correct trajectories. Experimental results show that our agent has superior navigation performance on Landmark-RxR, en-RxR and R2R. Our dataset and code are available at <https://github.com/hekj/Landmark-RxR>.

A Winning Hand: Compressing Deep Networks Can Improve Out-of-Distribution Robustness

- James Diffenderfer, Brian Bartoldson, Shreya Chaganti, Jize Zhang, Bhavya Kailkhura
- abstract@[open-review](#): Successful adoption of deep learning (DL) in the wild requires models to be: (1) compact, (2) accurate, and (3) robust to distributional shifts. Unfortunately, efforts towards simultaneously meeting these requirements have mostly been unsuccessful. This raises an important question: Is the inability to create Compact, Accurate, and Robust Deep neural networks (CARDs) fundamental? To answer this question, we perform a large-scale analysis of popular model compression techniques which uncovers several intriguing patterns. Notably, in contrast to traditional pruning approaches (e.g., fine tuning and gradual magnitude pruning), we find that ``lottery ticket-style'' approaches can surprisingly be used to produce CARDs, including binary-weight CARDs. Specifically, we are able to create extremely compact CARDs that, compared to their larger counterparts, have similar test accuracy and matching (or better) robustness---simply by pruning and (optionally) quantizing. Leveraging the compactness of CARDs, we develop a simple domain-adaptive test-time ensembling approach (CARD-Decks) that uses a gating module to dynamically select appropriate CARDs from the CARD-Deck based on their spectral-similarity with test samples. The proposed approach builds a "winning hand" of CARDs that establishes a new state-of-the-art (on RobustBench) on CIFAR-10-C accuracies (i.e., 96.8% standard and 92.75% robust) and CIFAR-100-C accuracies (80.6% standard and 71.3% robust) with better memory usage than non-compressed baselines (pretrained CARDs and CARD-Decks available at <https://github.com/RobustBench/robustbench>). Finally, we provide theoretical support for our empirical findings.

On the Importance of Gradients for Detecting Distributional Shifts in the Wild

- Rui Huang, Andrew Geng, Yixuan Li
- abstract@[open-review](#): Detecting out-of-distribution (OOD) data has become a critical component in ensuring the safe deployment of machine learning models in the real world. Existing OOD detection approaches primarily rely on the output or feature space for deriving OOD scores, while largely overlooking information from the gradient space. In this paper, we present GradNorm, a simple and effective approach for detecting OOD inputs by utilizing information extracted from the gradient space. GradNorm directly employs the vector norm of gradients, backpropagated from the KL divergence

between the softmax output and a uniform probability distribution. Our key idea is that the magnitude of gradients is higher for in-distribution (ID) data than that for OOD data, making it informative for OOD detection. GradNorm demonstrates superior performance, reducing the average FPR95 by up to 16.33% compared to the previous best method.

[Iterative Methods for Private Synthetic Data: Unifying Framework and New Methods](#)

- Terrance Liu, Giuseppe Vietri, Steven Z. Wu
- abstract@[open-review](#): We study private synthetic data generation for query release, where the goal is to construct a sanitized version of a sensitive dataset, subject to differential privacy, that approximately preserves the answers to a large collection of statistical queries. We first present an algorithmic framework that unifies a long line of iterative algorithms in the literature. Under this framework, we propose two new methods. The first method, private entropy projection (PEP), can be viewed as an advanced variant of MWEM that adaptively reuses past query measurements to boost accuracy. Our second method, generative networks with the exponential mechanism (GEM), circumvents computational bottlenecks in algorithms such as MWEM and PEP by optimizing over generative models parameterized by neural networks, which capture a rich family of distributions while enabling fast gradient-based optimization. We demonstrate that PEP and GEM empirically outperform existing algorithms. Furthermore, we show that GEM nicely incorporates prior information from public data while overcoming limitations of PMW^APub, the existing state-of-the-art method that also leverages public data.

[Understanding End-to-End Model-Based Reinforcement Learning Methods as Implicit Parameterization](#)

- Clement Gehring, Kenji Kawaguchi, Jiaoyang Huang, Leslie Kaelbling
- abstract@[open-review](#): Estimating the per-state expected cumulative rewards is a critical aspect of reinforcement learning approaches, however the experience is obtained, but standard deep neural-network function-approximation methods are often inefficient in this setting. An alternative approach, exemplified by value iteration networks, is to learn transition and reward models of a latent Markov decision process whose value predictions fit the data. This approach has been shown empirically to converge faster to a more robust solution in many cases, but there has been little theoretical study of this phenomenon. In this paper, we explore such implicit representations of value functions via theory and focused experimentation. We prove that, for a linear parametrization, gradient descent converges to global optima despite non-linearity and non-convexity introduced by the implicit representation. Furthermore, we derive convergence rates for both cases which allow us to identify conditions under which stochastic gradient descent (SGD) with this implicit representation converges substantially faster than its explicit counterpart. Finally, we provide empirical results in some simple domains that illustrate the theoretical findings.

[Mirror Langevin Monte Carlo: the Case Under Isoperimetry](#)

- Qijia Jiang
- abstract@[open-review](#): Motivated by the connection between sampling and optimization, we study a mirror descent analogue of Langevin dynamics and analyze three different discretization schemes, giving nonasymptotic convergence rate under functional inequalities such as Log-Sobolev in the corresponding metric. Compared to the Euclidean setting, the result reveals intricate relationship between the underlying geometry and the target distribution and suggests that care might need to be taken in order for the discretized algorithm to achieve vanishing bias with diminishing stepsize for sampling from potentials under weaker smoothness/convexity regularity conditions.

[Do Different Tracking Tasks Require Different Appearance Models?](#)

- Zhongdao Wang, Hengshuang Zhao, Ya-Li Li, Shengjin Wang, Philip Torr, Luca Bertinetto
- abstract@[open-review](#): Tracking objects of interest in a video is one of the most popular and widely applicable problems in computer vision. However, with the years, a Cambrian explosion of use cases and benchmarks has fragmented the problem in a multitude of different experimental setups. As a consequence, the literature has fragmented too, and now novel approaches proposed by the community are usually specialised to fit only one specific setup. To understand to what extent this specialisation is necessary, in this work we present UniTrack, a solution to address five different tasks within the same framework. UniTrack consists of a single and task-agnostic appearance model, which can be learned in a supervised or self-supervised fashion, and multiple ``heads'' that address individual tasks and do not require training. We show how most tracking tasks can be solved within this framework, and that the same appearance model can be successfully used to obtain results that are competitive against specialised methods for most of the tasks considered. The framework also allows us to analyse appearance models obtained with the most recent self-supervised methods, thus extending their evaluation and comparison to a larger variety of important problems.

[Towards robust vision by multi-task learning on monkey visual cortex](#)

- Shahd Safarani, Arne Nix, Konstantin Willeke, Santiago Cadena, Kelli Restivo, George Denfield, Andreas Tolias, Fabian Sinz
- abstract@[open-review](#): Deep neural networks set the state-of-the-art across many tasks in computer vision, but their generalization ability to simple image distortions is surprisingly fragile. In contrast, the mammalian visual system is robust to a wide range of perturbations. Recent work suggests that this generalization ability can be explained by useful inductive biases encoded in the representations of visual stimuli throughout the visual cortex. Here, we successfully leveraged these inductive biases with a multi-task learning approach: we jointly trained a deep network to perform image classification and to predict neural activity in macaque primary visual cortex (V1) in response to the same natural stimuli. We measured the out-of-distribution generalization abilities of our resulting network by testing its robustness to common image distortions. We found that co-training on monkey V1 data indeed leads to increased robustness despite the absence of those distortions during training. Additionally, we showed that our network's robustness is often very close to that of an Oracle network where parts of the architecture are directly trained on noisy images. Our results also demonstrated that the network's representations become more brain-like as their robustness improves. Using a novel constrained reconstruction analysis, we investigated what makes our brain-regularized network more robust. We found that our monkey co-trained network is more sensitive to content than noise when compared to a Baseline network that we trained for image classification alone. Using DeepGaze-predicted saliency maps for ImageNet images, we found that the monkey co-trained network tends to be more sensitive to salient regions in a scene, reminiscent of existing theories on the role of V1 in the detection of object borders and bottom-up saliency. Overall, our work expands the promising research avenue of transferring inductive biases from biological to artificial neural networks on the representational level, and provides a novel analysis of the effects of our transfer.

[Arbitrary Conditional Distributions with Energy](#)

- Ryan Strauss, Junier B. Oliva
- abstract@[open-review](#): Modeling distributions of covariates, or density estimation, is a core challenge in unsupervised learning. However, the majority of work only considers the joint distribution, which has limited relevance to practical situations. A more general and useful problem is arbitrary conditional density estimation, which aims to model any possible conditional distribution over a set of covariates, reflecting the more realistic setting of inference based on prior knowledge. We propose a novel method, Arbitrary Conditioning with Energy (ACE), that can simultaneously estimate the distribution $\Pr(\mathbf{x}_u \mid \mathbf{x}_o)$ for all possible subsets of unobserved features \mathbf{x}_u and observed features \mathbf{x}_o . ACE is designed to avoid unnecessary bias and complexity --- we specify densities with a highly expressive energy function and reduce the problem to only learning one-dimensional conditionals (from which more complex distributions can be recovered during inference). This results in an approach that is both simpler and higher-performing than prior methods. We show that ACE achieves state-of-the-art for arbitrary conditional likelihood estimation and data imputation on standard benchmarks.

[Learning Domain Invariant Representations in Goal-conditioned Block MDPs](#)

- Beining Han, Chongyi Zheng, Harris Chan, Keiran Paster, Michael Zhang, Jimmy Ba
- abstract@[open-review](#): Deep Reinforcement Learning (RL) is successful in solving many complex Markov Decision Processes (MDPs) problems. However, agents often face unanticipated environmental changes after deployment in the real world. These changes are often spurious and unrelated to the underlying problem, such as background shifts for visual input agents. Unfortunately, deep RL policies are usually sensitive to these changes and fail to act robustly against them. This resembles the problem of domain generalization in supervised learning. In this work, we study this problem for goal-conditioned RL agents. We propose a theoretical framework in the Block MDP setting that characterizes the generalizability of goal-conditioned policies to new environments. Under this framework, we develop a practical method PA-SkewFit that enhances domain generalization. The empirical evaluation shows that our goal-conditioned RL agent can perform well in various unseen test environments, improving by 50% over baselines.

[Near-Optimal Multi-Perturbation Experimental Design for Causal Structure Learning](#)

- Scott Sussex, Caroline Uhler, Andreas Krause
- abstract@[open-review](#): Causal structure learning is a key problem in many domains. Causal structures can be learnt by performing experiments on the system of interest. We address the largely unexplored problem of designing a batch of experiments that each simultaneously intervene on multiple variables. While potentially more informative than the commonly considered single-variable interventions, selecting such interventions is algorithmically much more challenging, due to the doubly-exponential combinatorial search space over sets of composite interventions. In this paper, we develop efficient algorithms for optimizing different objective functions quantifying the informativeness of a budget-constrained batch of experiments. By establishing novel submodularity properties of these objectives, we provide approximation guarantees for our algorithms. Our algorithms empirically perform superior to both random interventions and algorithms that only select single-variable interventions.

[Fuzzy Clustering with Similarity Queries](#)

- Wasim Huleihel, Arya Mazumdar, Soumyabrata Pal
- abstract@[open-review](#): The fuzzy or soft k -means objective is a popular generalization of the well-known k -means problem, extending the clustering capability of the k -means to datasets that are uncertain, vague and otherwise hard to cluster. In this paper, we propose a semi-supervised active clustering framework, where the learner is allowed to interact with an oracle (domain expert), asking for the similarity between a certain set of chosen items. We study the query and computational complexities of clustering in this framework. We prove that having a few of such similarity queries enables one to get a polynomial-time approximation algorithm to an otherwise conjecturally NP-hard problem. In particular, we provide algorithms for fuzzy clustering in this setting that ask $O(\mathsf{poly}(k)\log n)$ similarity queries and run with polynomial-time-complexity, where n is the number of items. The fuzzy k -means objective is nonconvex, with k -means as a special case, and is equivalent to some other generic nonconvex problem such as non-negative matrix factorization. The ubiquitous Lloyd-type algorithms (or alternating-minimization algorithms) can get stuck at a local minima. Our results show that by making few similarity queries, the problem becomes easier to solve. Finally, we test our algorithms over real-world datasets, showing their effectiveness in real-world applications.

[Improving black-box optimization in VAE latent space using decoder uncertainty](#)

- Pascal Notin, José Miguel Hernández-Lobato, Yarin Gal
- abstract@[open-review](#): Optimization in the latent space of variational autoencoders is a promising approach to generate high-dimensional discrete objects that maximize an expensive black-box property (e.g., drug-likeness in molecular generation, function approximation with arithmetic expressions). However, existing methods lack robustness as they may decide to explore areas of the latent space for which no data was available during training and where the decoder can be unreliable, leading to the generation of unrealistic or invalid objects. We propose to leverage the epistemic uncertainty of the decoder to guide the optimization process. This is not trivial though, as a naive estimation of uncertainty in the high-dimensional and structured settings we consider would result in high estimator variance. To solve this problem, we introduce an importance sampling-based estimator that provides more robust estimates of epistemic uncertainty. Our uncertainty-guided optimization approach does not require modifications of the model architecture nor the training process. It produces samples with a better trade-off between black-box objective and validity of the generated samples, sometimes improving both simultaneously. We illustrate these advantages across several experimental settings in digit generation, arithmetic expression approximation and molecule generation for drug design.

[Sample Selection for Fair and Robust Training](#)

- Yuji Roh, Kangwook Lee, Steven Whang, Changho Suh
- abstract@[open-review](#): Fairness and robustness are critical elements of Trustworthy AI that need to be addressed together. Fairness is about learning an unbiased model while robustness is about learning from corrupted data, and it is known that addressing only one of them may have an adverse affect on the other. In this work, we propose a sample selection-based algorithm for fair and robust training. To this end, we formulate a combinatorial optimization problem for the unbiased selection of samples in the presence of data corruption. Observing that solving this optimization problem is strongly NP-hard, we propose a greedy algorithm that is efficient and effective in practice. Experiments show that our method obtains fairness and robustness that are better than or comparable to the state-of-the-art technique, both on synthetic and benchmark real datasets. Moreover, unlike other fair and robust training baselines, our algorithm can be used by only modifying the sampling step in batch selection without changing the training algorithm or leveraging additional clean data.

[NeurWIN: Neural Whittle Index Network For Restless Bandits Via Deep RL](#)

- Khaled Nakhleh, Santosh Ganji, Ping-Chun Hsieh, I-Hong Hou, Srinivas Shakkottai
- abstract@[open-review](#): Whittle index policy is a powerful tool to obtain asymptotically optimal solutions for the notoriously intractable problem of restless bandits. However, finding the Whittle indices remains a difficult problem for many practical restless bandits with convoluted transition kernels. This paper proposes NeurWIN, a neural Whittle index network that seeks to learn the Whittle indices for any restless bandits by leveraging mathematical properties of the Whittle indices. We show that a neural network that produces the Whittle index is also one that produces the optimal control for a set of Markov decision problems. This property motivates using deep reinforcement learning for the training of NeurWIN. We demonstrate the utility of NeurWIN by evaluating its performance for three recently studied restless bandit problems. Our experiment results show that the performance of NeurWIN is significantly better than other RL algorithms.

[Sageflow: Robust Federated Learning against Both Stragglers and Adversaries](#)

- Jungwuk Park, Dong-Jun Han, Minseok Choi, Jaekyun Moon
- abstract@[open-review](#): While federated learning (FL) allows efficient model training with local data at edge devices, among major issues still to be resolved are: slow devices known as stragglers and malicious attacks launched by adversaries. While the presence of both of these issues raises serious concerns in practical FL systems, no known schemes or combinations of schemes effectively address them at the same time. We propose Sageflow, staleness-aware grouping with entropy-based filtering and loss-weighted averaging, to handle both stragglers and adversaries simultaneously. Model grouping and weighting according to staleness (arrival delay) provides robustness against stragglers, while entropy-based filtering and loss-weighted

averaging, working in a highly complementary fashion at each grouping stage, counter a wide range of adversary attacks. A theoretical bound is established to provide key insights into the convergence behavior of Sageflow. Extensive experimental results show that Sageflow outperforms various existing methods aiming to handle stragglers/adversaries.

[Alias-Free Generative Adversarial Networks](#)

- Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, Timo Aila
- abstract@[open-review](#): We observe that despite their hierarchical convolutional nature, the synthesis process of typical generative adversarial networks depends on absolute pixel coordinates in an unhealthy manner. This manifests itself as, e.g., detail appearing to be glued to image coordinates instead of the surfaces of depicted objects. We trace the root cause to careless signal processing that causes aliasing in the generator network. Interpreting all signals in the network as continuous, we derive generally applicable, small architectural changes that guarantee that unwanted information cannot leak into the hierarchical synthesis process. The resulting networks match the FID of StyleGAN2 but differ dramatically in their internal representations, and they are fully equivariant to translation and rotation even at subpixel scales. Our results pave the way for generative models better suited for video and animation.

[Noise2Score: Tweedie's Approach to Self-Supervised Image Denoising without Clean Images](#)

- Kwanyoung Kim, Jong Chul Ye
- abstract@[open-review](#): Recently, there has been extensive research interest in training deep networks to denoise images without clean reference. However, the representative approaches such as Noise2Noise, Noise2Void, Stein's unbiased risk estimator (SURE), etc. seem to differ from one another and it is difficult to find the coherent mathematical structure. To address this, here we present a novel approach, called Noise2Score, which reveals a missing link in order to unite these seemingly different approaches. Specifically, we show that image denoising problems without clean images can be addressed by finding the mode of the posterior distribution and that the Tweedie's formula offers an explicit solution through the score function (i.e. the gradient of loglikelihood). Our method then uses the recent finding that the score function can be stably estimated from the noisy images using the amortized residual denoising autoencoder, the method of which is closely related to Noise2Noise or Nose2Void. Our Noise2Score approach is so universal that the same network training can be used to remove noises from images that are corrupted by any exponential family distributions and noise parameters. Using extensive experiments with Gaussian, Poisson, and Gamma noises, we show that Noise2Score significantly outperforms the state-of-the-art self-supervised denoising methods in the benchmark data set such as (C)BSD68, Set12, and Kodak, etc.

[Continuous Mean-Covariance Bandits](#)

- Yihan Du, Siwei Wang, Zhixuan Fang, Longbo Huang
- abstract@[open-review](#): Existing risk-aware multi-armed bandit models typically focus on risk measures of individual options such as variance. As a result, they cannot be directly applied to important real-world online decision making problems with correlated options. In this paper, we propose a novel Continuous Mean-Covariance Bandit (CMCB) model to explicitly take into account option correlation. Specifically, in CMCB, there is a learner who sequentially chooses weight vectors on given options and observes random feedback according to the decisions. The agent's objective is to achieve the best trade-off between reward and risk, measured with option covariance. To capture different reward observation scenarios in practice, we consider three feedback settings, i.e., full-information, semi-bandit and full-bandit feedback. We propose novel algorithms with optimal regrets (within logarithmic factors), and provide matching lower bounds to validate their optimality. The experimental results also demonstrate the superiority of our algorithms. To the best of our knowledge, this is the first work that considers option correlation in risk-aware bandits and explicitly quantifies how arbitrary covariance structures impact the learning performance. The novel analytical techniques we developed for exploiting the estimated covariance to build concentration and bounding the risk of selected actions based on sampling strategy properties can likely find applications in other bandit analysis and be of independent interests.

[Dynamic Visual Reasoning by Learning Differentiable Physics Models from Video and Language](#)

- Mingyu Ding, Zhenfang Chen, Tao Du, Ping Luo, Josh Tenenbaum, Chuang Gan
- abstract@[open-review](#): In this work, we propose a unified framework, called Visual Reasoning with Differentiable Physics (VRDP), that can jointly learn visual concepts and infer physics models of objects and their interactions from videos and language. This is achieved by seamlessly integrating three components: a visual perception module, a concept learner, and a differentiable physics engine. The visual perception module parses each video frame into object-centric trajectories and represents them as latent scene representations. The concept learner grounds visual concepts (e.g., color, shape, and material) from these object-centric representations based on the language, thus providing prior knowledge for the physics engine. The differentiable physics model, implemented as an impulse-based differentiable rigid-body simulator, performs differentiable physical simulation based on the grounded concepts to infer physical properties, such as mass, restitution, and velocity, by fitting the simulated trajectories into the video observations. Consequently, these learned concepts and physical models can explain what we have seen and imagine what is about to happen in future and counterfactual scenarios. Integrating differentiable physics into the dynamic reasoning framework offers several appealing benefits. More accurate dynamics prediction in learned physics models enables state-of-the-art performance on both synthetic and real-world benchmarks while still maintaining high transparency and interpretability; most notably, VRDP improves the accuracy of predictive and counterfactual questions by 4.5% and 11.5% compared to its best counterpart. VRDP is also highly data-efficient: physical parameters can be optimized from very few videos, and even a single video can be sufficient. Finally, with all physical parameters inferred, VRDP can quickly learn new concepts from a few examples.

[Solving Soft Clustering Ensemble via \\$k\\$-Sparse Discrete Wasserstein Barycenter](#)

- Ruizhe Qin, Mengying Li, Hu Ding
- abstract@[open-review](#): Clustering ensemble is one of the most important problems in ensemble learning. Though it has been extensively studied in the past decades, the existing methods often suffer from the issues like high computational complexity and the difficulty on understanding the consensus. In this paper, we study the more general soft clustering ensemble problem where each individual solution is a soft clustering. We connect it to the well-known discrete Wasserstein barycenter problem in geometry. Based on some novel geometric insights in high dimensions, we propose the sampling-based algorithms with provable quality guarantees. We also provide the systematical analysis on the consensus of our model. Finally, we conduct the experiments to evaluate our proposed algorithms.

[Bayesian Adaptation for Covariate Shift](#)

- Aurick Zhou, Sergey Levine
- abstract@[open-review](#): When faced with distribution shift at test time, deep neural networks often make inaccurate predictions with unreliable uncertainty estimates. While improving the robustness of neural networks is one promising approach to mitigate this issue, an appealing alternate to robustifying networks against all possible test-time shifts is to instead directly adapt them to unlabeled inputs from the particular distribution shift we encounter at test time. However, this poses a challenging question: in the standard Bayesian model for supervised learning, unlabeled inputs are conditionally independent of model parameters when the labels are unobserved, so what can unlabeled data tell us about the model parameters at test-time? In this paper, we derive a Bayesian model that provides for a well-defined relationship between unlabeled inputs under distributional shift and model parameters, and show how approximate inference in this model can be instantiated with a simple regularized entropy minimization procedure at test-time. We evaluate our method on a variety of distribution shifts for image classification, including image corruptions, natural distribution shifts, and domain adaptation settings, and show that our method improves both accuracy and uncertainty estimation.

Perturb-and-max-product: Sampling and learning in discrete energy-based models

- Miguel Lazaro-Gredilla, Antoine Dedieu, Dileep George
- abstract@[open-review](#): Perturb-and-MAP offers an elegant approach to approximately sample from a energy-based model (EBM) by computing the maximum-a-posteriori (MAP) configuration of a perturbed version of the model. Sampling in turn enables learning. However, this line of research has been hindered by the general intractability of the MAP computation. Very few works venture outside tractable models, and when they do, they use linear programming approaches, which as we will show, have several limitations. In this work we present perturb-and-max-product (PMP), a parallel and scalable mechanism for sampling and learning in discrete EBMs. Models can be arbitrary as long as they are built using tractable factors. We show that (a) for Ising models, PMP is orders of magnitude faster than Gibbs and Gibbs-with-Gradients (GWG) at learning and generating samples of similar or better quality; (b) PMP is able to learn and sample from RBMs; (c) in a large, entangled graphical model in which Gibbs and GWG fail to mix, PMP succeeds.

Towards Unifying Behavioral and Response Diversity for Open-ended Learning in Zero-sum Games

- Xiangyu Liu, Hangtian Jia, Ying Wen, Yujing Hu, Yingfeng Chen, Changjie Fan, ZHIPENG HU, Yaodong Yang
- abstract@[open-review](#): Measuring and promoting policy diversity is critical for solving games with strong non-transitive dynamics where strategic cycles exist, and there is no consistent winner (e.g., Rock-Paper-Scissors). With that in mind, maintaining a pool of diverse policies via open-ended learning is an attractive solution, which can generate auto-curricula to avoid being exploited. However, in conventional open-ended learning algorithms, there are no widely accepted definitions for diversity, making it hard to construct and evaluate the diverse policies. In this work, we summarize previous concepts of diversity and work towards offering a unified measure of diversity in multi-agent open-ended learning to include all elements in Markov games, based on both Behavioral Diversity (BD) and Response Diversity (RD). At the trajectory distribution level, we re-define BD in the state-action space as the discrepancies of occupancy measures. For the reward dynamics, we propose RD to characterize diversity through the responses of policies when encountering different opponents. We also show that many current diversity measures fall in one of the categories of BD or RD but not both. With this unified diversity measure, we design the corresponding diversity-promoting objective and population effectivity when seeking the best responses in open-ended learning. We validate our methods in both relatively simple games like matrix game, non-transitive mixture model, and the complex \textit{Google Research Football} environment. The population found by our methods reveals the lowest exploitability, highest population effectivity in matrix game and non-transitive mixture model, as well as the largest goal difference when interacting with opponents of various levels in \textit{Google Research Football}.

Towards Better Understanding of Training Certifiably Robust Models against Adversarial Examples

- Sungyoon Lee, Woojin Lee, Jinseong Park, Jaewook Lee
- abstract@[open-review](#): We study the problem of training certifiably robust models against adversarial examples. Certifiable training minimizes an upper bound on the worst-case loss over the allowed perturbation, and thus the tightness of the upper bound is an important factor in building certifiably robust models. However, many studies have shown that Interval Bound Propagation (IBP) training uses much looser bounds but outperforms other models that use tighter bounds. We identify another key factor that influences the performance of certifiable training: \textit{smoothness of the loss landscape}. We find significant differences in the loss landscapes across many linear relaxation-based methods, and that the current state-of-the-arts method often has a landscape with favorable optimization properties. Moreover, to test the claim, we design a new certifiable training method with the desired properties. With the tightness and the smoothness, the proposed method achieves a decent performance under a wide range of perturbations, while others with only one of the two factors can perform well only for a specific range of perturbations. Our code is available at \url{https://github.com/sungyoon-lee/LossLandscapeMatters}.

Mitigating Covariate Shift in Imitation Learning via Offline Data With Partial Coverage

- Jonathan Chang, Masatoshi Uehara, Dhruv Sreenivas, Rahul Kidambi, Wen Sun
- abstract@[open-review](#): This paper studies offline Imitation Learning (IL) where an agent learns to imitate an expert demonstrator without additional online environment interactions. Instead, the learner is presented with a static offline dataset of state-action-next state triples from a potentially less proficient behavior policy. We introduce Model-based IL from Offline data (MILO): an algorithmic framework that utilizes the static dataset to solve the offline IL problem efficiently both in theory and in practice. In theory, even if the behavior policy is highly sub-optimal compared to the expert, we show that as long as the data from the behavior policy provides sufficient coverage on the expert state-action traces (and with no necessity for a global coverage over the entire state-action space), MILO can provably combat the covariate shift issue in IL. Complementing our theory results, we also demonstrate that a practical implementation of our approach mitigates covariate shift on benchmark MuJoCo continuous control tasks. We demonstrate that with behavior policies whose performances are less than half of that of the expert, MILO still successfully imitates with an extremely low number of expert state-action pairs while traditional offline IL methods such as behavior cloning (BC) fail completely. Source code is provided at <https://github.com/jdchang1/milo>.

Global Filter Networks for Image Classification

- Yongming Rao, Wenliang Zhao, Zheng Zhu, Jiwen Lu, Jie Zhou
- abstract@[open-review](#): Recent advances in self-attention and pure multi-layer perceptrons (MLP) models for vision have shown great potential in achieving promising performance with fewer inductive biases. These models are generally based on learning interaction among spatial locations from raw data. The complexity of self-attention and MLP grows quadratically as the image size increases, which makes these models hard to scale up when high-resolution features are required. In this paper, we present the Global Filter Network (GFNet), a conceptually simple yet computationally efficient architecture, that learns long-term spatial dependencies in the frequency domain with log-linear complexity. Our architecture replaces the self-attention layer in vision transformers with three key operations: a 2D discrete Fourier transform, an element-wise multiplication between frequency-domain features and learnable global filters, and a 2D inverse Fourier transform. We exhibit favorable accuracy/complexity trade-offs of our models on both ImageNet and downstream tasks. Our results demonstrate that GFNet can be a very competitive alternative to transformer-style models and CNNs in efficiency, generalization ability and robustness. Code is available at <https://github.com/raoyongming/GFNet>

CAFE: Catastrophic Data Leakage in Vertical Federated Learning

- Xiao Jin, Pin-Yu Chen, Chia-Yi Hsu, Chia-Mu Yu, Tianyi Chen
- abstract@[open-review](#): Recent studies show that private training data can be leaked through the gradients sharing mechanism deployed in distributed machine learning systems, such as federated learning (FL). Increasing batch size to complicate data recovery is often viewed as a promising defense strategy against data leakage. In this paper, we revisit this defense premise and propose an advanced data leakage attack with theoretical justification to efficiently recover batch data from the shared aggregated gradients. We name our proposed method as catastrophic data leakage in vertical federated learning (CAFE). Comparing to existing data leakage attacks, our extensive experimental results on vertical FL settings demonstrate the effectiveness of CAFE to perform large-batch data leakage attack with improved data recovery quality. We also propose a practical countermeasure to mitigate CAFE. Our results suggest that private data participated in standard FL, especially the vertical case, have a high risk of being leaked from the training gradients. Our analysis implies unprecedented and practical data leakage risks in those learning settings. The code of our work is available at <https://github.com/DeRafael/CAFE>.

Fault-Tolerant Federated Reinforcement Learning with Theoretical Guarantee

- Xiaofeng Fan, Yining Ma, Zhongxiang Dai, Wei Jing, Cheston Tan, Bryan Kian Hsiang Low
- abstract@[open-review](#): The growing literature of Federated Learning (FL) has recently inspired Federated Reinforcement Learning (FRL) to encourage multiple agents to federatively build a better decision-making policy without sharing raw trajectories. Despite its promising applications, existing works on FRL fail to I) provide theoretical analysis on its convergence, and II) account for random system failures and adversarial attacks. Towards this end, we propose the first FRL framework the convergence of which is guaranteed and tolerant to less than half of the participating agents being random system failures or adversarial attackers. We prove that the sample efficiency of the proposed framework is guaranteed to improve with the number of agents and is able to account for such potential failures or attacks. All theoretical results are empirically verified on various RL benchmark tasks.

[Compacter: Efficient Low-Rank Hypercomplex Adapter Layers](#)

- Rabeeh Karimi Mahabadi, James Henderson, Sebastian Ruder
- abstract@[open-review](#): Adapting large-scale pretrained language models to downstream tasks via fine-tuning is the standard method for achieving state-of-the-art performance on NLP benchmarks. However, fine-tuning all weights of models with millions or billions of parameters is sample-inefficient, unstable in low-resource settings, and wasteful as it requires storing a separate copy of the model for each task. Recent work has developed parameter-efficient fine-tuning methods, but these approaches either still require a relatively large number of parameters or underperform standard fine-tuning. In this work, we propose Compacter, a method for fine-tuning large-scale language models with a better trade-off between task performance and the number of trainable parameters than prior work. Compacter accomplishes this by building on top of ideas from adapters, low-rank optimization, and parameterized hypercomplex multiplication layers. Specifically, Compacter inserts task-specific weight matrices into a pretrained model's weights, which are computed efficiently as a sum of Kronecker products between shared "slow" weights and fast rank-one matrices defined per Compacter layer. By only training 0.047% of a pretrained model's parameters, Compacter performs on par with standard fine-tuning on GLUE and outperforms standard fine-tuning on SuperGLUE and low-resource settings. Our code is publicly available at <https://github.com/rabeehk/compacter>.

[Distilling Image Classifiers in Object Detectors](#)

- Shuxuan Guo, Jose M. Alvarez, Mathieu Salzmann
- abstract@[open-review](#): Knowledge distillation constitutes a simple yet effective way to improve the performance of a compact student network by exploiting the knowledge of a more powerful teacher. Nevertheless, the knowledge distillation literature remains limited to the scenario where the student and the teacher tackle the same task. Here, we investigate the problem of transferring knowledge not only across architectures but also across tasks. To this end, we study the case of object detection and, instead of following the standard detector-to-detector distillation approach, introduce a classifier-to-detector knowledge transfer framework. In particular, we propose strategies to exploit the classification teacher to improve both the detector's recognition accuracy and localization performance. Our experiments on several detectors with different backbones demonstrate the effectiveness of our approach, allowing us to outperform the state-of-the-art detector-to-detector distillation methods.

[Subgroup Generalization and Fairness of Graph Neural Networks](#)

- Jiaqi Ma, Junwei Deng, Qiaozhu Mei
- abstract@[open-review](#): Despite enormous successful applications of graph neural networks (GNNs), theoretical understanding of their generalization ability, especially for node-level tasks where data are not independent and identically-distributed (IID), has been sparse. The theoretical investigation of the generalization performance is beneficial for understanding fundamental issues (such as fairness) of GNN models and designing better learning methods. In this paper, we present a novel PAC-Bayesian analysis for GNNs under a non-IID semi-supervised learning setup. Moreover, we analyze the generalization performances on different subgroups of unlabeled nodes, which allows us to further study an accuracy-(dis)parity-style (un)fairness of GNNs from a theoretical perspective. Under reasonable assumptions, we demonstrate that the distance between a test subgroup and the training set can be a key factor affecting the GNN performance on that subgroup, which calls special attention to the training node selection for fair learning. Experiments across multiple GNN models and datasets support our theoretical results.

[Scaling Neural Tangent Kernels via Sketching and Random Features](#)

- Amir Zandieh, Insu Han, Haim Avron, Neta Shoham, Chaewon Kim, Jinwoo Shin
- abstract@[open-review](#): The Neural Tangent Kernel (NTK) characterizes the behavior of infinitely-wide neural networks trained under least squares loss by gradient descent. Recent works also report that NTK regression can outperform finitely-wide neural networks trained on small-scale datasets. However, the computational complexity of kernel methods has limited its use in large-scale learning tasks. To accelerate learning with NTK, we design a near input-sparsity time approximation algorithm for NTK, by sketching the polynomial expansions of arc-cosine kernels: our sketch for the convolutional counterpart of NTK (CNTK) can transform any image using a linear runtime in the number of pixels. Furthermore, we prove a spectral approximation guarantee for the NTK matrix, by combining random features (based on leverage score sampling) of the arc-cosine kernels with a sketching algorithm. We benchmark our methods on various large-scale regression and classification tasks and show that a linear regressor trained on our CNTK features matches the accuracy of exact CNTK on CIFAR-10 dataset while achieving 150x speedup.

[BatchQuant: Quantized-for-all Architecture Search with Robust Quantizer](#)

- Haoping Bai, Meng Cao, Ping Huang, Jiulong Shan
- abstract@[open-review](#): As the applications of deep learning models on edge devices increase at an accelerating pace, fast adaptation to various scenarios with varying resource constraints has become a crucial aspect of model deployment. As a result, model optimization strategies with adaptive configuration are becoming increasingly popular. While single-shot quantized neural architecture search enjoys flexibility in both model architecture and quantization policy, the combined search space comes with many challenges, including instability when training the weight-sharing supernet and difficulty in navigating the exponentially growing search space. Existing methods tend to either limit the architecture search space to a small set of options or limit the quantization policy search space to fixed precision policies. To this end, we propose BatchQuant, a robust quantizer formulation that allows fast and stable training of a compact, single-shot, mixed-precision, weight-sharing supernet. We employ BatchQuant to train a compact supernet (offering over \$10^{76}\$ quantized subnets) within substantially fewer GPU hours than previous methods. Our approach, Quantized-for-all (QFA), is the first to seamlessly extend one-shot weight-sharing NAS supernet to support subnets with arbitrary ultra-low bitwidth mixed-precision quantization policies without retraining. QFA opens up new possibilities in joint hardware-aware neural architecture search and quantization. We demonstrate the effectiveness of our method on ImageNet and achieve SOTA Top-1 accuracy under a low complexity constraint (<20 MFLOPs).

[Long Short-Term Transformer for Online Action Detection](#)

- Mingze Xu, Yuanjun Xiong, Hao Chen, Xinyu Li, Wei Xia, Zhuowen Tu, Stefano Soatto
- abstract@[open-review](#): We present Long Short-term TRansformer (LSTR), a temporal modeling algorithm for online action detection, which employs a long- and short-term memory mechanism to model prolonged sequence data. It consists of an LSTR encoder that dynamically leverages coarse-scale historical information from an extended temporal window (e.g., 2048 frames spanning of up to 8 minutes), together with an LSTR decoder that focuses on a short time window (e.g., 32 frames spanning 8 seconds) to model the fine-scale characteristics of the data. Compared to prior work, LSTR provides an effective and efficient method to model long videos with fewer heuristics, which is validated by extensive empirical analysis. LSTR achieves state-of-the-art performance on three standard online action detection benchmarks, THUMOS'14, TVSeries, and HACS Segment. Code has been made available at: <https://xumingze0308.github.io/projects/lstr>.

Near Optimal Policy Optimization via REPS

- Aldo Pacchiano, Jonathan N Lee, Peter Bartlett, Ofir Nachum
- abstract@[open-review](#): Since its introduction a decade ago, relative entropy policy search (REPS) has demonstrated successful policy learning on a number of simulated and real-world robotic domains, not to mention providing algorithmic components used by many recently proposed reinforcement learning (RL) algorithms. While REPS is commonly known in the community, there exist no guarantees on its performance when using stochastic and gradient-based solvers. In this paper we aim to fill this gap by providing guarantees and convergence rates for the sub-optimality of a policy learned using first-order optimization methods applied to the REPS objective. We first consider the setting in which we are given access to exact gradients and demonstrate how near-optimality of the objective translates to near-optimality of the policy. We then consider the practical setting of stochastic gradients, and introduce a technique that uses generative access to the underlying Markov decision process to compute parameter updates that maintain favorable convergence to the optimal regularized policy.

Self-Consistent Models and Values

- Greg Farquhar, Kate Baumli, Zita Marinho, Angelos Filos, Matteo Hessel, Hado P. van Hasselt, David Silver
- abstract@[open-review](#): Learned models of the environment provide reinforcement learning (RL) agents with flexible ways of making predictions about the environment. Models enable planning, i.e. using more computation to improve value functions or policies, without requiring additional environment interactions. In this work, we investigate a way of augmenting model-based RL, by additionally encouraging a learned model and value function to be jointly self-consistent. This lies in contrast to classic planning methods like Dyna, which only update the value function to be consistent with the model. We propose a number of possible self-consistency updates, study them empirically in both the tabular and function approximation settings, and find that with appropriate choices self-consistency can be useful both for policy evaluation and control.

Learning on Random Balls is Sufficient for Estimating (Some) Graph Parameters

- Takanori Maehara, Hoang NT
- abstract@[open-review](#): Theoretical analyses for graph learning methods often assume a complete observation of the input graph. Such an assumption might not be useful for handling any-size graphs due to the scalability issues in practice. In this work, we develop a theoretical framework for graph classification problems in the partial observation setting (i.e., subgraph samplings). Equipped with insights from graph limit theory, we propose a new graph classification model that works on a randomly sampled subgraph and a novel topology to characterize the representability of the model. Our theoretical framework contributes a theoretical validation of mini-batch learning on graphs and leads to new learning-theoretic results on generalization bounds as well as size-generalizability without assumptions on the input.

Risk-Averse Bayes-Adaptive Reinforcement Learning

- Marc Rigter, Bruno Lacerda, Nick Hawes
- abstract@[open-review](#): In this work, we address risk-averse Bayes-adaptive reinforcement learning. We pose the problem of optimising the conditional value at risk (CVaR) of the total return in Bayes-adaptive Markov decision processes (MDPs). We show that a policy optimising CVaR in this setting is risk-averse to both the epistemic uncertainty due to the prior distribution over MDPs, and the aleatoric uncertainty due to the inherent stochasticity of MDPs. We reformulate the problem as a two-player stochastic game and propose an approximate algorithm based on Monte Carlo tree search and Bayesian optimisation. Our experiments demonstrate that our approach significantly outperforms baseline approaches for this problem.

Iterative Connecting Probability Estimation for Networks

- Yichen Qin, Linhan Yu, Yang Li
- abstract@[open-review](#): Estimating the probabilities of connections between vertices in a random network using an observed adjacency matrix is an important task for network data analysis. Many existing estimation methods are based on certain assumptions on network structure, which limit their applicability in practice. Without making strong assumptions, we develop an iterative connecting probability estimation method based on neighborhood averaging. Starting at a random initial point or an existing estimate, our method iteratively updates the pairwise vertex distances, the sets of similar vertices, and connecting probabilities to improve the precision of the estimate. We propose a two-stage neighborhood selection procedure to achieve the trade-off between smoothness of the estimate and the ability to discover local structure. The tuning parameters can be selected by cross-validation. We establish desirable theoretical properties for our method, and further justify its superior performance by comparing with existing methods in simulation and real data analysis.

Learning to Adapt via Latent Domains for Adaptive Semantic Segmentation

- Yunan Liu, Shanshan Zhang, Yang Li, Jian Yang
- abstract@[open-review](#): Domain adaptive semantic segmentation aims to transfer knowledge learned from labeled source domain to unlabeled target domain. To narrow down the domain gap and ease adaptation difficulty, some recent methods translate source images to target-like images (latent domains), which are used as supplement or substitute to the original source data. Nevertheless, these methods neglect to explicitly model the relationship of knowledge transferring across different domains. Alternatively, in this work we break through the standard “source-target” one pair adaptation framework and construct multiple adaptation pairs (e.g. “source-latent” and “latent-target”). The purpose is to use the meta-knowledge (how to adapt) learned from one pair as guidance to assist the adaptation of another pair under a meta-learning framework. Furthermore, we extend our method to a more practical setting of open compound domain adaptation (a.k.a multiple-target domain adaptation), where the target is a compound of multiple domains without domain labels. In this setting, we embed an additional pair of “latent-latent” to reduce the domain gap between the source and different latent domains, allowing the model to adapt well on multiple target domains simultaneously. When evaluated on standard benchmarks, our method is superior to the state-of-the-art methods in both the single target and multiple-target domain adaptation settings.

Single Layer Predictive Normalized Maximum Likelihood for Out-of-Distribution Detection

- Koby Bibas, Meir Feder, Tal Hassner
- abstract@[open-review](#): Detecting out-of-distribution (OOD) samples is vital for developing machine learning based models for critical safety systems. Common approaches for OOD detection assume access to some OOD samples during training which may not be available in a real-life scenario. Instead, we utilize the predictive normalized maximum likelihood (pNML) learner, in which no assumptions are made on the tested input. We derive an explicit expression of the pNML and its generalization error, denoted as the regret, for a single layer neural network (NN). We show that this learner generalizes well when (i) the test vector resides in a subspace spanned by the eigenvectors associated with the large eigenvalues of the empirical correlation matrix of the training data, or (ii) the test sample is far from the decision boundary. Furthermore, we describe how to efficiently apply the derived pNML regret to any pretrained deep NN, by employing the explicit pNML for the last layer, followed by the softmax function. Applying the derived regret to deep NN requires neither additional tunable parameters nor extra data. We extensively evaluate our approach on 74 OOD detection benchmarks using DenseNet-100, ResNet-34, and WideResNet-40 models trained with CIFAR-100, CIFAR-10, SVHN, and ImageNet-30 showing a significant improvement of up to 15.6% over recent leading methods.

[Prototypical Cross-Attention Networks for Multiple Object Tracking and Segmentation](#)

- Lei Ke, Xia Li, Martin Danelljan, Yu-Wing Tai, Chi-Keung Tang, Fisher Yu
- abstract@[open-review](#): Multiple object tracking and segmentation requires detecting, tracking, and segmenting objects belonging to a set of given classes. Most approaches only exploit the temporal dimension to address the association problem, while relying on single frame predictions for the segmentation mask itself. We propose Prototypical Cross-Attention Network (PCAN), capable of leveraging rich spatio-temporal information for online multiple object tracking and segmentation. PCAN first distills a space-time memory into a set of prototypes and then employs cross-attention to retrieve rich information from the past frames. To segment each object, PCAN adopts a prototypical appearance module to learn a set of contrastive foreground and background prototypes, which are then propagated over time. Extensive experiments demonstrate that PCAN outperforms current video instance tracking and segmentation competition winners on both Youtube-VIS and BDD100K datasets, and shows efficacy to both one-stage and two-stage segmentation frameworks. Code and video resources are available at <http://vis.xyz/pub/pcan>.

[Algorithmic Instabilities of Accelerated Gradient Descent](#)

- Amit Attia, Tomer Koren
- abstract@[open-review](#): We study the algorithmic stability of Nesterov's accelerated gradient method. For convex quadratic objectives, Chen et al. (2018) proved that the uniform stability of the method grows quadratically with the number of optimization steps, and conjectured that the same is true for the general convex and smooth case. We disprove this conjecture and show, for two notions of algorithmic stability (including uniform stability), that the stability of Nesterov's accelerated method in fact deteriorates exponentially fast with the number of gradient steps. This stands in sharp contrast to the bounds in the quadratic case, but also to known results for non-accelerated gradient methods where stability typically grows linearly with the number of steps.

[Learning Optimal Predictive Checklists](#)

- Haoran Zhang, Quaid Morris, Berk Ustun, Marzyeh Ghassemi
- abstract@[open-review](#): Checklists are simple decision aids that are often used to promote safety and reliability in clinical applications. In this paper, we present a method to learn checklists for clinical decision support. We represent predictive checklists as discrete linear classifiers with binary features and unit weights. We then learn globally optimal predictive checklists from data by solving an integer programming problem. Our method allows users to customize checklists to obey complex constraints, including constraints to enforce group fairness and to binarize real-valued features at training time. In addition, it pairs models with an optimality gap that can inform model development and determine the feasibility of learning sufficiently accurate checklists on a given dataset. We pair our method with specialized techniques that speed up its ability to train a predictive checklist that performs well and has a small optimality gap. We benchmark the performance of our method on seven clinical classification problems, and demonstrate its practical benefits by training a short-form checklist for PTSD screening. Our results show that our method can fit simple predictive checklists that perform well and that can easily be customized to obey a rich class of custom constraints.

[Finite Sample Analysis of Average-Reward TD Learning and \\$Q\\$-Learning](#)

- Sheng Zhang, Zhe Zhang, Siva Theja Maguluri
- abstract@[open-review](#): The focus of this paper is on sample complexity guarantees of average-reward reinforcement learning algorithms, which are known to be more challenging to study than their discounted-reward counterparts. To the best of our knowledge, we provide the first known finite sample guarantees using both constant and diminishing step sizes of (i) average-reward TD(λ) with linear function approximation for policy evaluation and (ii) average-reward Q -learning in the tabular setting to find the optimal policy. A major challenge is that since the value functions are agnostic to an additive constant, the corresponding Bellman operators are no longer contraction mappings under any norm. We obtain the results for TD(λ) by working in an appropriately defined subspace that ensures uniqueness of the solution. For Q -learning, we exploit the span seminorm contractive property of the Bellman operator, and construct a novel Lyapunov function obtained by infimal convolution of a generalized Moreau envelope and the indicator function of a set.

[Generalization Bounds for Graph Embedding Using Negative Sampling: Linear vs Hyperbolic](#)

- Atsushi Suzuki, Atsushi Nitanda, jing wang, Linchuan Xu, Kenji Yamanishi, Marc Cavazza
- abstract@[open-review](#): Graph embedding, which represents real-world entities in a mathematical space, has enabled numerous applications such as analyzing natural languages, social networks, biochemical networks, and knowledge bases. It has been experimentally shown that graph embedding in hyperbolic space can represent hierarchical tree-like data more effectively than embedding in linear space, owing to hyperbolic space's exponential growth property. However, since the theoretical comparison has been limited to ideal noiseless settings, the potential for the hyperbolic space's property to worsen the generalization error for practical data has not been analyzed. In this paper, we provide a generalization error bound applicable for graph embedding both in linear and hyperbolic spaces under various negative sampling settings that appear in graph embedding. Our bound states that error is polynomial and exponential with respect to the embedding space's radius in linear and hyperbolic spaces, respectively, which implies that hyperbolic space's exponential growth property worsens the error. Using our bound, we clarify the data size condition on which graph embedding in hyperbolic space can represent a tree better than in Euclidean space by discussing the bias-variance trade-off. Our bound also shows that imbalanced data distribution, which often appears in graph embedding, can worsen the error.

[Gradient Starvation: A Learning Proclivity in Neural Networks](#)

- Mohammad Pezeshki, Oumar Kaba, Yoshua Bengio, Aaron C. Courville, Doina Precup, Guillaume Lajoie
- abstract@[open-review](#): We identify and formalize a fundamental gradient descent phenomenon resulting in a learning proclivity in over-parameterized neural networks. Gradient Starvation arises when cross-entropy loss is minimized by capturing only a subset of features relevant for the task, despite the presence of other predictive features that fail to be discovered. This work provides a theoretical explanation for the emergence of such feature imbalance in neural networks. Using tools from Dynamical Systems theory, we identify simple properties of learning dynamics during gradient descent that lead to this imbalance, and prove that such a situation can be expected given certain statistical structure in training data. Based on our proposed formalism, we develop guarantees for a novel regularization method aimed at decoupling feature learning dynamics, improving accuracy and robustness in cases hindered by gradient starvation. We illustrate our findings with simple and real-world out-of-distribution (OOD) generalization experiments.

[Offline Reinforcement Learning as One Big Sequence Modeling Problem](#)

- Michael Janner, Qiyang Li, Sergey Levine
- abstract@[open-review](#): Reinforcement learning (RL) is typically viewed as the problem of estimating single-step policies (for model-free RL) or single-step models (for model-based RL), leveraging the Markov property to factorize the problem in time. However, we can also view RL as a sequence modeling problem: predict a sequence of actions that leads to a sequence of high rewards. Viewed in this way, it is tempting to consider whether powerful, high-capacity sequence prediction models that work well in other supervised learning domains, such as natural-language processing, can also provide simple and effective solutions to the RL problem. To this end, we explore how RL can be reframed as "one big sequence modeling" problem, using state-of-the-art Transformer architectures to model distributions over sequences of states, actions, and rewards. Addressing RL as a sequence modeling problem

significantly simplifies a range of design decisions: we no longer require separate behavior policy constraints, as is common in prior work on offline model-free RL, and we no longer require ensembles or other epistemic uncertainty estimators, as is common in prior work on model-based RL. All of these roles are filled by the same Transformer sequence model. In our experiments, we demonstrate the flexibility of this approach across imitation learning, goal-conditioned RL, and offline RL.

[Optimality and Stability in Federated Learning: A Game-theoretic Approach](#)

- Kate Donahue, Jon Kleinberg
- abstract@[open-review](#): Federated learning is a distributed learning paradigm where multiple agents, each only with access to local data, jointly learn a global model. There has recently been an explosion of research aiming not only to improve the accuracy rates of federated learning, but also provide certain guarantees around social good properties such as total error. One branch of this research has taken a game-theoretic approach, and in particular, prior work has viewed federated learning as a hedonic game, where error-minimizing players arrange themselves into federating coalitions. This past work proves the existence of stable coalition partitions, but leaves open a wide range of questions, including how far from optimal these stable solutions are. In this work, we motivate and define a notion of optimality given by the average error rates among federating agents (players). First, we provide and prove the correctness of an efficient algorithm to calculate an optimal (error minimizing) arrangement of players. Next, we analyze the relationship between the stability and optimality of an arrangement. First, we show that for some regions of parameter space, all stable arrangements are optimal (Price of Anarchy equal to 1). However, we show this is not true for all settings: there exist examples of stable arrangements with higher cost than optimal (Price of Anarchy greater than 1). Finally, we give the first constant-factor bound on the performance gap between stability and optimality, proving that the total error of the worst stable solution can be no higher than 9 times the total error of an optimal solution (Price of Anarchy bound of 9).

[Understanding Deflation Process in Over-parametrized Tensor Decomposition](#)

- Rong Ge, Yunwei Ren, Xiang Wang, Mo Zhou
- abstract@[open-review](#): In this paper we study the training dynamics for gradient flow on over-parametrized tensor decomposition problems. Empirically, such training process often first fits larger components and then discovers smaller components, which is similar to a tensor deflation process that is commonly used in tensor decomposition algorithms. We prove that for orthogonally decomposable tensor, a slightly modified version of gradient flow would follow a tensor deflation process and recover all the tensor components. Our proof suggests that for orthogonal tensors, gradient flow dynamics works similarly as greedy low-rank learning in the matrix setting, which is a first step towards understanding the implicit regularization effect of over-parametrized models for low-rank tensors.

[Privately Learning Subspaces](#)

- Vikrant Singhal, Thomas Steinke
- abstract@[open-review](#): Private data analysis suffers a costly curse of dimensionality. However, the data often has an underlying low-dimensional structure. For example, when optimizing via gradient descent, the gradients often lie in or near a low-dimensional subspace. If that low-dimensional structure can be identified, then we can avoid paying (in terms of privacy or accuracy) for the high ambient dimension. We present differentially private algorithms that take input data sampled from a low-dimensional linear subspace (possibly with a small amount of error) and output that subspace (or an approximation to it). These algorithms can serve as a pre-processing step for other procedures.

[On the Value of Interaction and Function Approximation in Imitation Learning](#)

- Nived Rajaraman, Yanjun Han, Lin Yang, Jingbo Liu, Jiantao Jiao, Kannan Ramchandran
- abstract@[open-review](#): We study the statistical guarantees for the Imitation Learning (IL) problem in episodic MDPs. Rajaraman et al. (2020) show an information theoretic lower bound that in the worst case, a learner which can even actively query the expert policy suffers from a suboptimality growing quadratically in the length of the horizon, $\$H\$$. We study imitation learning under the $\$\mu\$$ -recoverability assumption of Ross et al. (2011) which assumes that the difference in the $\$Q\$$ -value under the expert policy across different actions in a state do not deviate beyond $\$\mu\$$ from the maximum. We show that the reduction proposed by Ross et al. (2010) is statistically optimal: the resulting algorithm upon interacting with the MDP for $\$N\$$ episodes results in a suboptimality bound of $\$widetilde{O} \left(\mu \mathcal{S} H / N \right)$ which we show is optimal up to log-factors. In contrast, we show that any algorithm which does not interact with the MDP and uses an offline dataset of $\$N\$$ expert trajectories must incur suboptimality growing as $\$gt�im \mathcal{S} H^2 / N\$$ even under the $\$\mu\$$ -recoverability assumption. This establishes a clear and provable separation of the minimax rates between the active setting and the no-interaction setting. We also study IL with linear function approximation. When the expert plays actions according to a linear classifier of known state-action features, we use the reduction to multi-class classification to show that with high probability, the suboptimality of behavior cloning is $\$widetilde{O}(dH^2/N)$ given $\$N\$$ rollouts from the optimal policy. This is optimal up to log-factors but can be improved to $\$widetilde{O}(dH/N)$ if we have a linear expert with parameter-sharing across time steps. In contrast, when the MDP transition structure is known to the learner such as in the case of simulators, we demonstrate fundamental differences compared to the tabular setting in terms of the performance of an optimal algorithm, Mimic-MD (Rajaraman et al. (2020)) when extended to the function approximation setting. Here, we introduce a new problem called confidence set linear classification, that can be used to construct sample-efficient IL algorithms.

[Shapeshifter: a Parameter-efficient Transformer using Factorized Reshaped Matrices](#)

- Aliakbar Panahi, Seyran Saeedi, Tom Arodz
- abstract@[open-review](#): Language models employ a very large number of trainable parameters. Despite being highly overparameterized, these networks often achieve good out-of-sample test performance on the original task and easily fine-tune to related tasks. Recent observations involving, for example, intrinsic dimension of the objective landscape and the lottery ticket hypothesis, indicate that often training actively involves only a small fraction of the parameter space. Thus, a question remains how large a parameter space needs to be in the first place — the evidence from recent work on model compression, parameter sharing, factorized representations, and knowledge distillation increasingly shows that models can be made much smaller and still perform well. Here, we focus on factorized representations of matrices that underpin dense, embedding, and self-attention layers. We use low-rank factorized representation of a reshaped and rearranged original matrix to achieve space efficient and expressive linear layers. We prove that stacking such low-rank layers increases their expressiveness, providing theoretical understanding for their effectiveness in deep networks. In Transformer models, our approach leads to more than ten-fold reduction in the number of total trainable parameters, including embedding, attention, and feed-forward layers, with little degradation in on-task performance. The approach operates out-of-the-box, replacing each parameter matrix with its compact equivalent while maintaining the architecture of the network.

[The Adaptive Doubly Robust Estimator and a Paradox Concerning Logging Policy](#)

- Masahiro Kato, Kenichiro McAlinn, Shota Yasui
- abstract@[open-review](#): The doubly robust (DR) estimator, which consists of two nuisance parameters, the conditional mean outcome and the logging policy (the probability of choosing an action), is crucial in causal inference. This paper proposes a DR estimator for dependent samples obtained from adaptive experiments. To obtain an asymptotically normal semiparametric estimator from dependent samples without non-Donsker nuisance estimators, we propose adaptive-fitting as a variant of sample-splitting. We also report an empirical paradox that our proposed DR estimator tends to show better performances compared to other estimators utilizing the true logging policy. While a similar phenomenon is known for estimators with i.i.d. samples,

traditional explanations based on asymptotic efficiency cannot elucidate our case with dependent samples. We confirm this hypothesis through simulation studies.

[Regularized Softmax Deep Multi-Agent Q-Learning](#)

- Ling Pan, Tabish Rashid, Bei Peng, Longbo Huang, Shimon Whiteson
- abstract@[open-review](#): Tackling overestimation in $\$Q\$$ -learning is an important problem that has been extensively studied in single-agent reinforcement learning, but has received comparatively little attention in the multi-agent setting. In this work, we empirically demonstrate that QMIX, a popular $\$Q\$$ -learning algorithm for cooperative multi-agent reinforcement learning (MARL), suffers from a more severe overestimation in practice than previously acknowledged, and is not mitigated by existing approaches. We rectify this with a novel regularization-based update scheme that penalizes large joint action-values that deviate from a baseline and demonstrate its effectiveness in stabilizing learning. Furthermore, we propose to employ a softmax operator, which we efficiently approximate in a novel way in the multi-agent setting, to further reduce the potential overestimation bias. Our approach, Regularized Softmax (RES) Deep Multi-Agent $\$Q\$$ -Learning, is general and can be applied to any $\$Q\$$ -learning based MARL algorithm. We demonstrate that, when applied to QMIX, RES avoids severe overestimation and significantly improves performance, yielding state-of-the-art results in a variety of cooperative multi-agent tasks, including the challenging StarCraft II micromanagement benchmarks.

[Physics-Aware Downsampling with Deep Learning for Scalable Flood Modeling](#)

- Niv Giladi, Zvika Ben-Haim, Sella Nevo, Yossi Matias, Daniel Soudry
- abstract@[open-review](#): Background. Floods are the most common natural disaster in the world, affecting the lives of hundreds of millions. Flood forecasting is therefore a vitally important endeavor, typically achieved using physical water flow simulations, which rely on accurate terrain elevation maps. However, such simulations, based on solving partial differential equations, are computationally prohibitive on a large scale. This scalability issue is commonly alleviated using a coarse grid representation of the elevation map, though this representation may distort crucial terrain details, leading to significant inaccuracies in the simulation.\Contributions. We train a deep neural network to perform physics-informed downsampling of the terrain map: we optimize the coarse grid representation of the terrain maps, so that the flood prediction will match the fine grid solution. For the learning process to succeed, we configure a dataset specifically for this task. We demonstrate that with this method, it is possible to achieve a significant reduction in computational cost, while maintaining an accurate solution. A reference implementation accompanies the paper as well as documentation and code for dataset reproduction.

[Systematic Generalization with Edge Transformers](#)

- Leon Bergen, Timothy O'Donnell, Dzmitry Bahdanau
- abstract@[open-review](#): Recent research suggests that systematic generalization in natural language understanding remains a challenge for state-of-the-art neural models such as Transformers and Graph Neural Networks. To tackle this challenge, we propose Edge Transformer, a new model that combines inspiration from Transformers and rule-based symbolic AI. The first key idea in Edge Transformers is to associate vector states with every edge, that is, with every pair of input nodes---as opposed to just every node, as it is done in the Transformer model. The second major innovation is a triangular attention mechanism that updates edge representations in a way that is inspired by unification from logic programming. We evaluate Edge Transformer on compositional generalization benchmarks in relational reasoning, semantic parsing, and dependency parsing. In all three settings, the Edge Transformer outperforms Relation-aware, Universal and classical Transformer baselines.

[TransformerFusion: Monocular RGB Scene Reconstruction using Transformers](#)

- Aljaz Bozic, Pablo Palafox, Justus Thies, Angela Dai, Matthias Niessner
- abstract@[open-review](#): We introduce TransformerFusion, a transformer-based 3D scene reconstruction approach. From an input monocular RGB video, the video frames are processed by a transformer network that fuses the observations into a volumetric feature grid representing the scene; this feature grid is then decoded into an implicit 3D scene representation. Key to our approach is the transformer architecture that enables the network to learn to attend to the most relevant image frames for each 3D location in the scene, supervised only by the scene reconstruction task. Features are fused in a coarse-to-fine fashion, storing fine-level features only where needed, requiring lower memory storage and enabling fusion at interactive rates. The feature grid is then decoded to a higher-resolution scene reconstruction, using an MLP-based surface occupancy prediction from interpolated coarse-to-fine 3D features. Our approach results in an accurate surface reconstruction, outperforming state-of-the-art multi-view stereo depth estimation methods, fully-convolutional 3D reconstruction approaches, and approaches using LSTM- or GRU-based recurrent networks for video sequence fusion.

[Maximum Likelihood Training of Score-Based Diffusion Models](#)

- Yang Song, Conor Durkan, Iain Murray, Stefano Ermon
- abstract@[open-review](#): Score-based diffusion models synthesize samples by reversing a stochastic process that diffuses data to noise, and are trained by minimizing a weighted combination of score matching losses. The log-likelihood of score-based diffusion models can be tractably computed through a connection to continuous normalizing flows, but log-likelihood is not directly optimized by the weighted combination of score matching losses. We show that for a specific weighting scheme, the objective upper bounds the negative log-likelihood, thus enabling approximate maximum likelihood training of score-based diffusion models. We empirically observe that maximum likelihood training consistently improves the likelihood of score-based diffusion models across multiple datasets, stochastic processes, and model architectures. Our best models achieve negative log-likelihoods of 2.83 and 3.76 bits/dim on CIFAR-10 and ImageNet \$32\times 32\$ without any data augmentation, on a par with state-of-the-art autoregressive models on these tasks.

[Global Convergence of Gradient Descent for Asymmetric Low-Rank Matrix Factorization](#)

- Tian Ye, Simon S. Du
- abstract@[open-review](#): We study the asymmetric low-rank factorization problem:
$$\min_{\mathbf{U}, \mathbf{V}} \|\mathbf{U}\|_{\mathbb{R}} \|\mathbf{V}\|_{\mathbb{R}} \text{ s.t. } \mathbf{U}^T \mathbf{V} = \mathbf{R}$$
where \mathbf{R} is a given matrix of size $m \times n$ and rank d . This is a canonical problem that admits two difficulties in optimization: 1) non-convexity and 2) non-smoothness (due to unbalancedness of \mathbf{U} and \mathbf{V}). This is also a prototype for more complex problems such as asymmetric matrix sensing and matrix completion. Despite being non-convex and non-smooth, it has been observed empirically that the randomly initialized gradient descent algorithm can solve this problem in polynomial time. Existing theories to explain this phenomenon all require artificial modifications of the algorithm, such as adding noise in each iteration and adding a balancing regularizer to balance the \mathbf{U} and \mathbf{V} . This paper presents the first proof that shows randomly initialized gradient descent converges to a global minimum of the asymmetric low-rank factorization problem with a polynomial rate. For the proof, we develop 1) a new symmetrization technique to capture the magnitudes of the symmetry and asymmetry, and 2) a quantitative perturbation analysis to approximate matrix derivatives. We believe both are useful for other related non-convex problems.

[Adaptive Data Augmentation on Temporal Graphs](#)

- Yiwei Wang, Yujun Cai, Yuxuan Liang, Henghui Ding, Changhu Wang, Siddharth Bhatia, Bryan Hooi

- abstract@[open-review](#): Temporal Graph Networks (TGNs) are powerful on modeling temporal graph data based on their increased complexity. Higher complexity carries with it a higher risk of overfitting, which makes TGNs capture random noise instead of essential semantic information. To address this issue, our idea is to transform the temporal graphs using data augmentation (DA) with adaptive magnitudes, so as to effectively augment the input features and preserve the essential semantic information. Based on this idea, we present the MeTA (Memory Tower Augmentation) module: a multi-level module that processes the augmented graphs of different magnitudes on separate levels, and performs message passing across levels to provide adaptively augmented inputs for every prediction. MeTA can be flexibly applied to the training of popular TGNs to improve their effectiveness without increasing their time complexity. To complement MeTA, we propose three DA strategies to realistically model noise by modifying both the temporal and topological features. Empirical results on standard datasets show that MeTA yields significant gains for the popular TGN models on edge prediction and node classification in an efficient manner.

[Regularized Frank-Wolfe for Dense CRFs: Generalizing Mean Field and Beyond](#)

- D.Khuê Lê-Huu, Karteek Alahari
- abstract@[open-review](#): We introduce regularized Frank-Wolfe, a general and effective algorithm for inference and learning of dense conditional random fields (CRFs). The algorithm optimizes a nonconvex continuous relaxation of the CRF inference problem using vanilla Frank-Wolfe with approximate updates, which are equivalent to minimizing a regularized energy function. Our proposed method is a generalization of existing algorithms such as mean field or concave-convex procedure. This perspective not only offers a unified analysis of these algorithms, but also allows an easy way of exploring different variants that potentially yield better performance. We illustrate this in our empirical results on standard semantic segmentation datasets, where several instantiations of our regularized Frank-Wolfe outperform mean field inference, both as a standalone component and as an end-to-end trainable layer in a neural network. We also show that dense CRFs, coupled with our new algorithms, produce significant improvements over strong CNN baselines.

[Terra: Imperative-Symbolic Co-Execution of Imperative Deep Learning Programs](#)

- Taebum Kim, Eunji Jeong, Geon-Woo Kim, Yunmo Koo, Sehoon Kim, Gyeongin Yu, Byung-Gon Chun
- abstract@[open-review](#): Imperative programming allows users to implement their deep neural networks (DNNs) easily and has become an essential part of recent deep learning (DL) frameworks. Recently, several systems have been proposed to combine the usability of imperative programming with the optimized performance of symbolic graph execution. Such systems convert imperative Python DL programs to optimized symbolic graphs and execute them. However, they cannot fully support the usability of imperative programming. For example, if an imperative DL program contains a Python feature with no corresponding symbolic representation (e.g., third-party library calls or unsupported dynamic control flows) they fail to execute the program. To overcome this limitation, we propose Terra, an imperative-symbolic co-execution system that can handle any imperative DL programs while achieving the optimized performance of symbolic graph execution. To achieve this, Terra builds a symbolic graph by decoupling DL operations from Python features. Then, Terra conducts the imperative execution to support all Python features, while delegating the decoupled operations to the symbolic execution. We evaluated Terra's performance improvement and coverage with ten imperative DL programs for several DNN architectures. The results show that Terra can speed up the execution of all ten imperative DL programs, whereas AutoGraph, one of the state-of-the-art systems, fails to execute five of them.

[Uniform Sampling over Episode Difficulty](#)

- Sébastien Arnold, Guneet Dhillon, Avinash Ravichandran, Stefano Soatto
- abstract@[open-review](#): Episodic training is a core ingredient of few-shot learning to train models on tasks with limited labelled data. Despite its success, episodic training remains largely understudied, prompting us to ask the question: what is the best way to sample episodes? In this paper, we first propose a method to approximate episode sampling distributions based on their difficulty. Building on this method, we perform an extensive analysis and find that sampling uniformly over episode difficulty outperforms other sampling schemes, including curriculum and easy-/hard-mining. As the proposed sampling method is algorithm agnostic, we can leverage these insights to improve few-shot learning accuracies across many episodic training algorithms. We demonstrate the efficacy of our method across popular few-shot learning datasets, algorithms, network architectures, and protocols.

[Scalable Intervention Target Estimation in Linear Models](#)

- Burak Varici, Karthikeyan Shanmugam, Prasanna Sattigeri, Ali Tajer
- abstract@[open-review](#): This paper considers the problem of estimating the unknown intervention targets in a causal directed acyclic graph from observational and interventional data. The focus is on soft interventions in linear structural equation models (SEMs). Current approaches to causal structure learning either work with known intervention targets or use hypothesis testing to discover the unknown intervention targets even for linear SEMs. This severely limits their scalability and sample complexity. This paper proposes a scalable and efficient algorithm that consistently identifies all intervention targets. The pivotal idea is to estimate the intervention sites from the difference between the precision matrices associated with the observational and interventional datasets. It involves repeatedly estimating such sites in different subsets of variables. The proposed algorithm can be used to also update a given observational Markov equivalence class into the interventional Markov equivalence class. Consistency, Markov equivalency, and sample complexity are established analytically. Finally, simulation results on both real and synthetic data demonstrate the gains of the proposed approach for scalable causal structure recovery. Implementation of the algorithm and the code to reproduce the simulation results are available at \url{https://github.com/bvarici/intervention-estimation}.

[Play to Grade: Testing Coding Games as Classifying Markov Decision Process](#)

- Allen Nie, Emma Brunskill, Chris Piech
- abstract@[open-review](#): Contemporary coding education often presents students with the task of developing programs that have user interaction and complex dynamic systems, such as mouse based games. While pedagogically compelling, there are no contemporary autonomous methods for providing feedback. Notably, interactive programs are impossible to grade by traditional unit tests. In this paper we formalize the challenge of providing feedback to interactive programs as a task of classifying Markov Decision Processes (MDPs). Each student's program fully specifies an MDP where the agent needs to operate and decide, under reasonable generalization, if the dynamics and reward model of the input MDP should be categorized as correct or broken. We demonstrate that by designing a cooperative objective between an agent and an autoregressive model, we can use the agent to sample differential trajectories from the input MDP that allows a classifier to determine membership: Play to Grade. Our method enables an automatic feedback system for interactive code assignments. We release a dataset of 711,274 anonymized student submissions to a single assignment with hand-coded bug labels to support future research.

[Distributional Reinforcement Learning for Multi-Dimensional Reward Functions](#)

- Pushi Zhang, Xiaoyu Chen, Li Zhao, Wei Xiong, Tao Qin, Tie-Yan Liu
- abstract@[open-review](#): A growing trend for value-based reinforcement learning (RL) algorithms is to capture more information than scalar value functions in the value network. One of the most well-known methods in this branch is distributional RL, which models return distribution instead of scalar value. In another line of work, hybrid reward architectures (HRA) in RL have studied to model source-specific value functions for each source of reward, which is also shown to be beneficial in performance. To fully inherit the benefits of distributional RL and hybrid reward architectures, we introduce Multi-Dimensional Distributional DQN (MD3QN), which extends distributional RL to model the joint return distribution from multiple reward sources. As a by-product of joint distribution modeling, MD3QN can capture not only the randomness in returns for each source of reward, but also the rich reward correlation between the randomness of different sources. We prove the convergence for the joint distributional Bellman operator and build our empirical

algorithm by minimizing the Maximum Mean Discrepancy between joint return distribution and its Bellman target. In experiments, our method accurately models the joint return distribution in environments with richly correlated reward functions, and outperforms previous RL methods utilizing multi-dimensional reward functions in the control setting.

[Differentiable Unsupervised Feature Selection based on a Gated Laplacian](#)

- Ofir Lindenbaum, Uri Shaham, Erez Peterfreund, Jonathan Svirsky, Nicolas Casey, Yuval Kluger
- abstract@[open-review](#): Scientific observations may consist of a large number of variables (features). Selecting a subset of meaningful features is often crucial for identifying patterns hidden in the ambient space. In this paper, we present a method for unsupervised feature selection, and we demonstrate its advantage in clustering, a common unsupervised task. We propose a differentiable loss that combines a graph Laplacian-based score that favors low-frequency features with a gating mechanism for removing nuisance features. Our method improves upon the naive graph Laplacian score by replacing it with a gated variant computed on a subset of low-frequency features. We identify this subset by learning the parameters of continuously relaxed Bernoulli variables, which gate the entire feature space. We mathematically motivate the proposed approach and demonstrate that it is crucial to compute the graph Laplacian on the gated inputs rather than on the full feature set in the high noise regime. Using several real-world examples, we demonstrate the efficacy and advantage of the proposed approach over leading baselines.

[Smooth Bilevel Programming for Sparse Regularization](#)

- Clarice Poon, Gabriel Peyré
- abstract@[open-review](#): Iteratively reweighted least square (IRLS) is a popular approach to solve sparsity-enforcing regression problems in machine learning. State of the art approaches are more efficient but typically rely on specific coordinate pruning schemes. In this work, we show how a surprisingly simple re-parametrization of IRLS, coupled with a bilevel resolution (instead of an alternating scheme) is able to achieve top performances on a wide range of sparsity (such as Lasso, group Lasso and trace norm regularizations), regularization strength (including hard constraints), and design matrices (ranging from correlated designs to differential operators). Similarly to IRLS, our method only involves linear systems resolutions, but in sharp contrast, corresponds to the minimization of a smooth function. Despite being non-convex, we show that there is no spurious minima and that saddle points are "ridable", so that there always exists a descent direction. We thus advocate for the use of a BFGS quasi-Newton solver, which makes our approach simple, robust and efficient. We perform a numerical benchmark of the convergence speed of our algorithm against state of the art solvers for Lasso, group Lasso, trace norm and linearly constrained problems. These results highlight the versatility of our approach, removing the need to use different solvers depending on the specificity of the ML problem under study.

[Grounding Representation Similarity Through Statistical Testing](#)

- Frances Ding, Jean-Stanislas Denain, Jacob Steinhardt
- abstract@[open-review](#): To understand neural network behavior, recent works quantitatively compare different networks' learned representations using canonical correlation analysis (CCA), centered kernel alignment (CKA), and other dissimilarity measures. Unfortunately, these widely used measures often disagree on fundamental observations, such as whether deep networks differing only in random initialization learn similar representations. These disagreements raise the question: which, if any, of these dissimilarity measures should we believe? We provide a framework to ground this question through a concrete test: measures should have \emph{sensitivity} to changes that affect functional behavior, and \emph{specificity} against changes that do not. We quantify this through a variety of functional behaviors including probing accuracy and robustness to distribution shift, and examine changes such as varying random initialization and deleting principal components. We find that current metrics exhibit different weaknesses, note that a classical baseline performs surprisingly well, and highlight settings where all metrics appear to fail, thus providing a challenge set for further improvement.

[A Consciousness-Inspired Planning Agent for Model-Based Reinforcement Learning](#)

- Mingde Zhao, Zhen Liu, Sitao Luan, Shuyuan Zhang, Doina Precup, Yoshua Bengio
- abstract@[open-review](#): We present an end-to-end, model-based deep reinforcement learning agent which dynamically attends to relevant parts of its state during planning. The agent uses a bottleneck mechanism over a set-based representation to force the number of entities to which the agent attends at each planning step to be small. In experiments, we investigate the bottleneck mechanism with several sets of customized environments featuring different challenges. We consistently observe that the design allows the planning agents to generalize their learned task-solving abilities in compatible unseen environments by attending to the relevant objects, leading to better out-of-distribution generalization performance.

[Reward-Free Model-Based Reinforcement Learning with Linear Function Approximation](#)

- Weitong ZHANG, Dongruo Zhou, Quanquan Gu
- abstract@[open-review](#): We study the model-based reward-free reinforcement learning with linear function approximation for episodic Markov decision processes (MDPs). In this setting, the agent works in two phases. In the exploration phase, the agent interacts with the environment and collects samples without the reward. In the planning phase, the agent is given a specific reward function and uses samples collected from the exploration phase to learn a good policy. We propose a new provably efficient algorithm, called UCRL-RFE under the Linear Mixture MDP assumption, where the transition probability kernel of the MDP can be parameterized by a linear function over certain feature mappings defined on the triplet of state, action, and next state. We show that to obtain an $\$epsilon$ -optimal policy for arbitrary reward function, UCRL-RFE needs to sample at most $\tilde{O}(H^5d^2\epsilon^{-2})$ episodes during the exploration phase. Here, H is the length of the episode, d is the dimension of the feature mapping. We also propose a variant of UCRL-RFE using Bernstein-type bonus and show that it needs to sample at most $\tilde{O}(H^4d(H+d)\epsilon^{-2})$ to achieve an $\$epsilon$ -optimal policy. By constructing a special class of linear Mixture MDPs, we also prove that for any reward-free algorithm, it needs to sample at least $\tilde{\Omega}(H^2d\epsilon^{-2})$ episodes to obtain an $\$epsilon$ -optimal policy. Our upper bound matches the lower bound in terms of the dependence on $\$epsilon$ and the dependence on d if $H \geq d$.

[Beltrami Flow and Neural Diffusion on Graphs](#)

- Benjamin Chamberlain, James Rowbottom, Davide Eynard, Francesco Di Giovanni, Xiaowen Dong, Michael Bronstein
- abstract@[open-review](#): We propose a novel class of graph neural networks based on the discretized Beltrami flow, a non-Euclidean diffusion PDE. In our model, node features are supplemented with positional encodings derived from the graph topology and jointly evolved by the Beltrami flow, producing simultaneously continuous feature learning, topology evolution. The resulting model generalizes many popular graph neural networks and achieves state-of-the-art results on several benchmarks.

[Think Big, Teach Small: Do Language Models Distil Occam's Razor?](#)

- Gonzalo Jaimovich-Lopez, David Castellano Falcón, Cesar Ferri, José Hernández-Orallo
- abstract@[open-review](#): Large language models have recently shown a remarkable ability for few-shot learning, including patterns of algorithmic nature. However, it is still an open question to determine what kind of patterns these models can capture and how many examples they need in their prompts. We frame this question as a teaching problem with strong priors, and study whether language models can identify simple algorithmic concepts from small witness sets. In particular, we explore how several GPT architectures, program induction systems and humans perform in terms of the complexity of the

concept and the number of additional examples, and how much their behaviour differs. This first joint analysis of language models and machine teaching can address key questions for artificial intelligence and machine learning, such as whether some strong priors, and Occam's razor in particular, can be distilled from data, making learning from a few examples possible.

[Disentangling Identifiable Features from Noisy Data with Structured Nonlinear ICA](#)

- Hermanni Hälvä, Sylvain Le Corff, Luc Lehéricy, Jonathan So, Yongjie Zhu, Elisabeth Gassiat, Aapo Hyvärinen
- abstract@[open-review](#): We introduce a new general identifiable framework for principled disentanglement referred to as Structured Nonlinear Independent Component Analysis (SNICA). Our contribution is to extend the identifiability theory of deep generative models for a very broad class of structured models. While previous works have shown identifiability for specific classes of time-series models, our theorems extend this to more general temporal structures as well as to models with more complex structures such as spatial dependencies. In particular, we establish the major result that identifiability for this framework holds even in the presence of noise of unknown distribution. Finally, as an example of our framework's flexibility, we introduce the first nonlinear ICA model for time-series that combines the following very useful properties: it accounts for both nonstationarity and autocorrelation in a fully unsupervised setting; performs dimensionality reduction; models hidden states; and enables principled estimation and inference by variational maximum-likelihood.

[Conditionally Parameterized, Discretization-Aware Neural Networks for Mesh-Based Modeling of Physical Systems](#)

- Jiayang Xu, Aniruddhe Pradhan, Karthikeyan Duraisamy
- abstract@[open-review](#): Simulations of complex physical systems are typically realized by discretizing partial differential equations (PDEs) on unstructured meshes. While neural networks have recently been explored for the surrogate and reduced order modeling of PDE solutions, they often ignore interactions or hierarchical relations between input features, and process them as concatenated mixtures. We generalize the idea of conditional parameterization -- using trainable functions of input parameters to generate the weights of a neural network, and extend them in a flexible way to encode critical information. Inspired by discretized numerical methods, choices of the parameters include physical quantities and mesh topology features. The functional relation between the modeled features and the parameters is built into the network architecture. The method is implemented on different networks and applied to frontier scientific machine learning tasks including the discovery of unmodeled physics, super-resolution of coarse fields, and the simulation of unsteady flows with chemical reactions. The results show that the conditionally-parameterized networks provide superior performance compared to their traditional counterparts. The CP-GNet - an architecture that can be trained on very few data snapshots - is proposed as the first deep learning model capable of standalone prediction of reacting flows on irregular meshes.

[USCO-Solver: Solving Undetermined Stochastic Combinatorial Optimization Problems](#)

- Guangmo Tong
- abstract@[open-review](#): Real-world decision-making systems are often subject to uncertainties that have to be resolved through observational data. Therefore, we are frequently confronted with combinatorial optimization problems of which the objective function is unknown and thus has to be debunked using empirical evidence. In contrast to the common practice that relies on a learning-and-optimization strategy, we consider the regression between combinatorial spaces, aiming to infer high-quality optimization solutions from samples of input-solution pairs -- without the need to learn the objective function. Our main deliverable is a universal solver that is able to handle abstract undetermined stochastic combinatorial optimization problems. For learning foundations, we present learning-error analysis under the PAC-Bayesian framework using a new margin-based analysis. In empirical studies, we demonstrate our design using proof-of-concept experiments, and compare it with other methods that are potentially applicable. Overall, we obtain highly encouraging experimental results for several classic combinatorial problems on both synthetic and real-world datasets.

[Adaptive Conformal Inference Under Distribution Shift](#)

- Isaac Gibbs, Emmanuel Candes
- abstract@[open-review](#): We develop methods for forming prediction sets in an online setting where the data generating distribution is allowed to vary over time in an unknown fashion. Our framework builds on ideas from conformal inference to provide a general wrapper that can be combined with any black box method that produces point predictions of the unseen label or estimated quantiles of its distribution. While previous conformal inference methods rely on the assumption that the data are exchangeable, our adaptive approach provably achieves the desired coverage frequency over long-time intervals irrespective of the true data generating process. We accomplish this by modelling the distribution shift as a learning problem in a single parameter whose optimal value is varying over time and must be continuously re-estimated. We test our method, adaptive conformal inference, on two real world datasets and find that its predictions are robust to visible and significant distribution shifts.

[Periodic Activation Functions Induce Stationarity](#)

- Lassi Meronen, Martin Trapp, Arno Solin
- abstract@[open-review](#): Neural network models are known to reinforce hidden data biases, making them unreliable and difficult to interpret. We seek to build models that 'know what they do not know' by introducing inductive biases in the function space. We show that periodic activation functions in Bayesian neural networks establish a connection between the prior on the network weights and translation-invariant, stationary Gaussian process priors. Furthermore, we show that this link goes beyond sinusoidal (Fourier) activations by also covering triangular wave and periodic ReLU activation functions. In a series of experiments, we show that periodic activation functions obtain comparable performance for in-domain data and capture sensitivity to perturbed inputs in deep neural networks for out-of-domain detection.

[Towards Optimal Strategies for Training Self-Driving Perception Models in Simulation](#)

- David Acuna, Jonah Philion, Sanja Fidler
- abstract@[open-review](#): Autonomous driving relies on a huge volume of real-world data to be labeled to high precision. Alternative solutions seek to exploit driving simulators that can generate large amounts of labeled data with a plethora of content variations. However, the domain gap between the synthetic and real data remains, raising the following important question: What are the best way to utilize a self-driving simulator for perception tasks?. In this work, we build on top of recent advances in domain-adaptation theory, and from this perspective, propose ways to minimize the reality gap. We primarily focus on the use of labels in the synthetic domain alone. Our approach introduces both a principled way to learn neural-invariant representations and a theoretically inspired view on how to sample the data from the simulator. Our method is easy to implement in practice as it is agnostic of the network architecture and the choice of the simulator. We showcase our approach on the bird's-eye-view vehicle segmentation task with multi-sensor data (cameras, lidar) using an open-source simulator (CARLA), and evaluate the entire framework on a real-world dataset (nuScenes). Last but not least, we show what types of variations (e.g. weather conditions, number of assets, map design and color diversity) matter to perception networks when trained with driving simulators, and which ones can be compensated for with our domain adaptation technique.

[KS-GNN: Keywords Search over Incomplete Graphs via Graphs Neural Network](#)

- YU HAO, Xin Cao, Yufan Sheng, Yixiang Fang, Wei Wang

- abstract@[open-review](#): Keyword search is a fundamental task to retrieve information that is the most relevant to the query keywords. Keyword search over graphs aims to find subtrees or subgraphs containing all query keywords ranked according to some criteria. Existing studies all assume that the graphs have complete information. However, real-world graphs may contain some missing information (such as edges or keywords), thus making the problem much more challenging. To solve the problem of keyword search over incomplete graphs, we propose a novel model named KS-GNN based on the graph neural network and the auto-encoder. By considering the latent relationships and the frequency of different keywords, the proposed KS-GNN aims to alleviate the effect of missing information and is able to learn low-dimensional representative node embeddings that preserve both graph structure and keyword features. Our model can effectively answer keyword search queries with linear time complexity over incomplete graphs. The experiments on four real-world datasets show that our model consistently achieves better performance than state-of-the-art baseline methods in graphs having missing information.

[Reconstruction for Powerful Graph Representations](#)

- Leonardo Cotta, Christopher Morris, Bruno Ribeiro
- abstract@[open-review](#): Graph neural networks (GNNs) have limited expressive power, failing to represent many graph classes correctly. While more expressive graph representation learning (GRL) alternatives can distinguish some of these classes, they are significantly harder to implement, may not scale well, and have not been shown to outperform well-tuned GNNs in real-world tasks. Thus, devising simple, scalable, and expressive GRL architectures that also achieve real-world improvements remains an open challenge. In this work, we show the extent to which graph reconstruction---reconstructing a graph from its subgraphs---can mitigate the theoretical and practical problems currently faced by GRL architectures. First, we leverage graph reconstruction to build two new classes of expressive graph representations. Secondly, we show how graph reconstruction boosts the expressive power of any GNN architecture while being a (provably) powerful inductive bias for invariances to vertex removals. Empirically, we show how reconstruction can boost GNN's expressive power---while maintaining its invariance to permutations of the vertices---by solving seven graph property tasks not solvable by the original GNN. Further, we demonstrate how it boosts state-of-the-art GNN's performance across nine real-world benchmark datasets.

[Revealing and Protecting Labels in Distributed Training](#)

- Trung Dang, Om Thakkar, Swaroop Ramaswamy, Rajiv Mathews, Peter Chin, Françoise Beaufays
- abstract@[open-review](#): Distributed learning paradigms such as federated learning often involve transmission of model updates, or gradients, over a network, thereby avoiding transmission of private data. However, it is possible for sensitive information about the training data to be revealed from such gradients. Prior works have demonstrated that labels can be revealed analytically from the last layer of certain models (e.g., ResNet), or they can be reconstructed jointly with model inputs by using Gradients Matching [Zhu et al.] with additional knowledge about the current state of the model. In this work, we propose a method to discover the set of labels of training samples from only the gradient of the last layer and the id to label mapping. Our method is applicable to a wide variety of model architectures across multiple domains. We demonstrate the effectiveness of our method for model training in two domains - image classification, and automatic speech recognition. Furthermore, we show that existing reconstruction techniques improve their efficacy when used in conjunction with our method. Conversely, we demonstrate that gradient quantization and sparsification can significantly reduce the success of the attack.

[Solving Graph-based Public Goods Games with Tree Search and Imitation Learning](#)

- Victor-Alexandru Darvariu, Stephen Hailes, Mirco Musolesi
- abstract@[open-review](#): Public goods games represent insightful settings for studying incentives for individual agents to make contributions that, while costly for each of them, benefit the wider society. In this work, we adopt the perspective of a central planner with a global view of a network of self-interested agents and the goal of maximizing some desired property in the context of a best-shot public goods game. Existing algorithms for this known NP-complete problem find solutions that are sub-optimal and cannot optimize for criteria other than social welfare. In order to efficiently solve public goods games, our proposed method directly exploits the correspondence between equilibria and the Maximal Independent Set (mIS) structural property of graphs. In particular, we define a Markov Decision Process which incrementally generates an mIS, and adopt a planning method to search for equilibria, outperforming existing methods. Furthermore, we devise a graph imitation learning technique that uses demonstrations of the search to obtain a graph neural network parametrized policy which quickly generalizes to unseen game instances. Our evaluation results show that this policy is able to reach 99.5% of the performance of the planning method while being three orders of magnitude faster to evaluate on the largest graphs tested. The methods presented in this work can be applied to a large class of public goods games of potentially high societal impact and more broadly to other graph combinatorial optimization problems.

[Stochastic Optimization of Areas Under Precision-Recall Curves with Provable Convergence](#)

- Qi Qi, Youzhi Luo, Zhao Xu, Shuiwang Ji, Tianbao Yang
- abstract@[open-review](#): Areas under ROC (AUROC) and precision-recall curves (AUPRC) are common metrics for evaluating classification performance for imbalanced problems. Compared with AUROC, AUPRC is a more appropriate metric for highly imbalanced datasets. While stochastic optimization of AUROC has been studied extensively, principled stochastic optimization of AUPRC has been rarely explored. In this work, we propose a principled technical method to optimize AUPRC for deep learning. Our approach is based on maximizing the averaged precision (AP), which is an unbiased point estimator of AUPRC. We cast the objective into a sum of dependent compositional functions with inner functions dependent on random variables of the outer level. We propose efficient adaptive and non-adaptive stochastic algorithms named SOAP with provable convergence guarantee under mild conditions by leveraging recent advances in stochastic compositional optimization. Extensive experimental results on image and graph datasets demonstrate that our proposed method outperforms prior methods on imbalanced problems in terms of AUPRC. To the best of our knowledge, our work represents the first attempt to optimize AUPRC with provable convergence. The SOAP has been implemented in the libAUC library at <https://libauc.org/>.

[Transfer Learning of Graph Neural Networks with Ego-graph Information Maximization](#)

- Qi Zhu, Carl Yang, Yidan Xu, Haonan Wang, Chao Zhang, Jiawei Han
- abstract@[open-review](#): Graph neural networks (GNNs) have achieved superior performance in various applications, but training dedicated GNNs can be costly for large-scale graphs. Some recent work started to study the pre-training of GNNs. However, none of them provide theoretical insights into the design of their frameworks, or clear requirements and guarantees towards their transferability. In this work, we establish a theoretically grounded and practically useful framework for the transfer learning of GNNs. Firstly, we propose a novel view towards the essential graph information and advocate the capturing of it as the goal of transferable GNN training, which motivates the design of EGI (Ego-Graph Information maximization) to analytically achieve this goal. Secondly, when node features are structure-relevant, we conduct an analysis of EGI transferability regarding the difference between the local graph Laplacians of the source and target graphs. We conduct controlled synthetic experiments to directly justify our theoretical conclusions. Comprehensive experiments on two real-world network datasets show consistent results in the analyzed setting of direct-transferring, while those on large-scale knowledge graphs show promising results in the more practical setting of transferring with fine-tuning.

[You are caught stealing my winning lottery ticket! Making a lottery ticket claim its ownership](#)

- Xuxi Chen, Tianlong Chen, Zhenyu Zhang, Zhangyang Wang

- abstract@[open-review](#): Despite tremendous success in many application scenarios, the training and inference costs of using deep learning are also rapidly increasing over time. The lottery ticket hypothesis (LTH) emerges as a promising framework to leverage a special sparse subnetwork (i.e., \$textit{winning ticket}\$) instead of a full model for both training and inference, that can lower both costs without sacrificing the performance. The main resource bottleneck of LTH is however the extraordinary cost to find the sparse mask of the winning ticket. That makes the found winning ticket become a valuable asset to the owners, highlighting the necessity of protecting its copyright. Our setting adds a new dimension to the recently soaring interest in protecting against the intellectual property (IP) infringement of deep models and verifying their ownerships, since they take owners' massive/unique resources to develop or train. While existing methods explored encrypted weights or predictions, we investigate a unique way to leverage sparse topological information to perform \$textit{lottery verification}\$, by developing several graph-based signatures that can be embedded as credentials. By further combining trigger set-based methods, our proposal can work in both white-box and black-box verification scenarios. Through extensive experiments, we demonstrate the effectiveness of lottery verification in diverse models (ResNet-20, ResNet-18, ResNet-50) on CIFAR-10 and CIFAR-100. Specifically, our verification is shown to be robust to removal attacks such as model fine-tuning and pruning, as well as several ambiguity attacks. Our codes are available at <https://github.com/VITA-Group/NO-stealing-LTH>.

[Complexity Lower Bounds for Nonconvex-Strongly-Concave Min-Max Optimization](#)

- Haochuan Li, Yi Tian, Jingzhao Zhang, Ali Jadbabaie
- abstract@[open-review](#): We provide a first-order oracle complexity lower bound for finding stationary points of min-max optimization problems where the objective function is smooth, nonconvex in the minimization variable, and strongly concave in the maximization variable. We establish a lower bound of $\Omega(\sqrt{\kappa}\epsilon^{-2})$ for deterministic oracles, where ϵ defines the level of approximate stationarity and κ is the condition number. Our lower bound matches the best existing upper bound in the ϵ and κ dependence up to logarithmic factors. For stochastic oracles, we provide a lower bound of $\Omega(\sqrt{\kappa}\epsilon^{-2} + \kappa^{1/3}\epsilon^{-4})$. It suggests that there is a gap between the best existing upper bound $O(\kappa^3\epsilon^{-4})$ and our lower bound in the condition number dependence.

[Early-stopped neural networks are consistent](#)

- Ziwei Ji, Justin Li, Matus Telgarsky
- abstract@[open-review](#): This work studies the behavior of shallow ReLU networks trained with the logistic loss via gradient descent on binary classification data where the underlying data distribution is general, and the (optimal) Bayes risk is not necessarily zero. In this setting, it is shown that gradient descent with early stopping achieves population risk arbitrarily close to optimal in terms of not just logistic and misclassification losses, but also in terms of calibration, meaning the sigmoid mapping of its outputs approximates the true underlying conditional distribution arbitrarily finely. Moreover, the necessary iteration, sample, and architectural complexities of this analysis all scale naturally with a certain complexity measure of the true conditional model. Lastly, while it is not shown that early stopping is necessary, it is shown that any classifier satisfying a basic local interpolation property is inconsistent.

[NxMTransformer: Semi-Structured Sparsification for Natural Language Understanding via ADMM](#)

- Connor Holmes, Minjia Zhang, Yuxiong He, Bo Wu
- abstract@[open-review](#): Natural Language Processing (NLP) has recently achieved great success by using huge pre-trained Transformer networks. However, these models often contain hundreds of millions or even billions of parameters, bringing challenges to online deployment due to latency constraints. Recently, hardware manufacturers have introduced dedicated hardware for NxM sparsity to provide the flexibility of unstructured pruning with the runtime efficiency of structured approaches. NxM sparsity permits arbitrarily selecting M parameters to retain from a contiguous group of N in the dense representation. However, due to the extremely high complexity of pre-trained models, the standard sparse fine-tuning techniques often fail to generalize well on downstream tasks, which have limited data resources. To address such an issue in a principled manner, we introduce a new learning framework, called NxMTransformer, to induce NxM semi-structured sparsity on pretrained language models for natural language understanding to obtain better performance. In particular, we propose to formulate the NxM sparsity as a constrained optimization problem and use Alternating Direction Method of Multipliers (ADMM) to optimize the downstream tasks while taking the underlying hardware constraints into consideration. ADMM decomposes the NxM sparsification problem into two sub-problems that can be solved sequentially, generating sparsified Transformer networks that achieve high accuracy while being able to effectively execute on newly released hardware. We apply our approach to a wide range of NLP tasks, and our proposed method is able to achieve 1.7 points higher accuracy in GLUE score than current best practices. Moreover, we perform detailed analysis on our approach and shed light on how ADMM affects fine-tuning accuracy for downstream tasks. Finally, we illustrate how NxMTransformer achieves additional performance improvement with knowledge distillation based methods.

[Reliable Decisions with Threshold Calibration](#)

- Roshni Sahoo, Shengjia Zhao, Alyssa Chen, Stefano Ermon
- abstract@[open-review](#): Decision makers rely on probabilistic forecasts to predict the loss of different decision rules before deployment. When the forecasted probabilities match the true frequencies, predicted losses will be accurate. Although perfect forecasts are typically impossible, probabilities can be calibrated to match the true frequencies on average. However, we find that this average notion of calibration, which is typically used in practice, does not necessarily guarantee accurate decision loss prediction. Specifically in the regression setting, the loss of threshold decisions, which are decisions based on whether the forecasted outcome falls above or below a cutoff, might not be predicted accurately. We propose a stronger notion of calibration called threshold calibration, which is exactly the condition required to ensure that decision loss is predicted accurately for threshold decisions. We provide an efficient algorithm which takes an uncalibrated forecaster as input and provably outputs a threshold-calibrated forecaster. Our procedure allows downstream decision makers to confidently estimate the loss of any threshold decision under any threshold loss function. Empirically, threshold calibration improves decision loss prediction without compromising on the quality of the decisions in two real-world settings: hospital scheduling decisions and resource allocation decisions.

[End-to-End Weak Supervision](#)

- Salva Rühling Cachay, Benedikt Boecking, Artur Dubrawski
- abstract@[open-review](#): Aggregating multiple sources of weak supervision (WS) can ease the data-labeling bottleneck prevalent in many machine learning applications, by replacing the tedious manual collection of ground truth labels. Current state of the art approaches that do not use any labeled training data, however, require two separate modeling steps: Learning a probabilistic latent variable model based on the WS sources -- making assumptions that rarely hold in practice -- followed by downstream model training. Importantly, the first step of modeling does not consider the performance of the downstream model. To address these caveats we propose an end-to-end approach for directly learning the downstream model by maximizing its agreement with probabilistic labels generated by reparameterizing previous probabilistic posteriors with a neural network. Our results show improved performance over prior work in terms of end model performance on downstream test sets, as well as in terms of improved robustness to dependencies among weak supervision sources.

[Shift Invariance Can Reduce Adversarial Robustness](#)

- Vasu Singla, Songwei Ge, Basri Ronen, David Jacobs

- abstract@[open-review](#): Shift invariance is a critical property of CNNs that improves performance on classification. However, we show that invariance to circular shifts can also lead to greater sensitivity to adversarial attacks. We first characterize the margin between classes when a shift-invariant { em linear} classifier is used. We show that the margin can only depend on the DC component of the signals. Then, using results about infinitely wide networks, we show that in some simple cases, fully connected and shift-invariant neural networks produce linear decision boundaries. Using this, we prove that shift invariance in neural networks produces adversarial examples for the simple case of two classes, each consisting of a single image with a black or white dot on a gray background. This is more than a curiosity; we show empirically that with real datasets and realistic architectures, shift invariance reduces adversarial robustness. Finally, we describe initial experiments using synthetic data to probe the source of this connection.

[Wisdom of the Crowd Voting: Truthful Aggregation of Voter Information and Preferences](#)

- Grant Schoenebeck, Biao Shuai Tao
- abstract@[open-review](#): We consider two-alternative elections where voters' preferences depend on a state variable that is not directly observable. Each voter receives a private signal that is correlated to the state variable. As a special case, our model captures the common scenario where voters can be categorized into three types: those who always prefer one alternative, those who always prefer the other, and those contingent voters whose preferences depends on the state. In this setting, even if every voter is a contingent voter, agents voting according to their private information need not result in the adoption of the universally preferred alternative, because the signals can be systematically biased. We present a mechanism that elicits and aggregates the private signals from the voters, and outputs the alternative that is favored by the majority. In particular, voters truthfully reporting their signals forms a strong Bayes Nash equilibrium (where no coalition of voters can deviate and receive a better outcome).

[Replay-Guided Adversarial Environment Design](#)

- Minqi Jiang, Michael Dennis, Jack Parker-Holder, Jakob Foerster, Edward Grefenstette, Tim Rocktäschel
- abstract@[open-review](#): Deep reinforcement learning (RL) agents may successfully generalize to new settings if trained on an appropriately diverse set of environment and task configurations. Unsupervised Environment Design (UED) is a promising self-supervised RL paradigm, wherein the free parameters of an underspecified environment are automatically adapted during training to the agent's capabilities, leading to the emergence of diverse training environments. Here, we cast Prioritized Level Replay (PLR), an empirically successful but theoretically unmotivated method that selectively samples randomly-generated training levels, as UED. We argue that by curating completely random levels, PLR, too, can generate novel and complex levels for effective training. This insight reveals a natural class of UED methods we call Dual Curriculum Design (DCD). Crucially, DCD includes both PLR and a popular UED algorithm, PAIRED, as special cases and inherits similar theoretical guarantees. This connection allows us to develop novel theory for PLR, providing a version with a robustness guarantee at Nash equilibria. Furthermore, our theory suggests a highly counterintuitive improvement to PLR: by stopping the agent from updating its policy on uncurated levels (training on less data), we can improve the convergence to Nash equilibria. Indeed, our experiments confirm that our new method, PLR\$^{\perp}\$, obtains better results on a suite of out-of-distribution, zero-shot transfer tasks, in addition to demonstrating that PLR\$^{\perp}\$ improves the performance of PAIRED, from which it inherited its theoretical framework.

[There Is No Turning Back: A Self-Supervised Approach for Reversibility-Aware Reinforcement Learning](#)

- Nathan Grinsztajn, Johan Ferret, Olivier Pietquin, Philippe Preux, Matthieu Geist
- abstract@[open-review](#): We propose to learn to distinguish reversible from irreversible actions for better informed decision-making in Reinforcement Learning (RL). From theoretical considerations, we show that approximate reversibility can be learned through a simple surrogate task: ranking randomly sampled trajectory events in chronological order. Intuitively, pairs of events that are always observed in the same order are likely to be separated by an irreversible sequence of actions. Conveniently, learning the temporal order of events can be done in a fully self-supervised way, which we use to estimate the reversibility of actions from experience, without any priors. We propose two different strategies that incorporate reversibility in RL agents, one strategy for exploration (RAE) and one strategy for control (RAC). We demonstrate the potential of reversibility-aware agents in several environments, including the challenging Sokoban game. In synthetic tasks, we show that we can learn control policies that never fail and reduce to zero the side-effects of interactions, even without access to the reward function.

[Learning to Execute: Efficient Learning of Universal Plan-Conditioned Policies in Robotics](#)

- Ingmar Schubert, Danny Driess, Ozgur S. Oguz, Marc Toussaint
- abstract@[open-review](#): Applications of Reinforcement Learning (RL) in robotics are often limited by high data demand. On the other hand, approximate models are readily available in many robotics scenarios, making model-based approaches like planning a data-efficient alternative. Still, the performance of these methods suffers if the model is imprecise or wrong. In this sense, the respective strengths and weaknesses of RL and model-based planners are complementary. In the present work, we investigate how both approaches can be integrated into one framework that combines their strengths. We introduce Learning to Execute (L2E), which leverages information contained in approximate plans to learn universal policies that are conditioned on plans. In our robotic manipulation experiments, L2E exhibits increased performance when compared to pure RL, pure planning, or baseline methods combining learning and planning.

[Self-Diagnosing GAN: Diagnosing Underrepresented Samples in Generative Adversarial Networks](#)

- Jinhee Lee, Haeri Kim, Youngkyu Hong, Hye Won Chung
- abstract@[open-review](#): Despite remarkable performance in producing realistic samples, Generative Adversarial Networks (GANs) often produce low-quality samples near low-density regions of the data manifold, e.g., samples of minor groups. Many techniques have been developed to improve the quality of generated samples, either by post-processing generated samples or by pre-processing the empirical data distribution, but at the cost of reduced diversity. To promote diversity in sample generation without degrading the overall quality, we propose a simple yet effective method to diagnose and emphasize underrepresented samples during training of a GAN. The main idea is to use the statistics of the discrepancy between the data distribution and the model distribution at each data instance. Based on the observation that the underrepresented samples have a high average discrepancy or high variability in discrepancy, we propose a method to emphasize those samples during training of a GAN. Our experimental results demonstrate that the proposed method improves GAN performance on various datasets, and it is especially effective in improving the quality and diversity of sample generation for minor groups.

[Online Multi-Armed Bandits with Adaptive Inference](#)

- Maria Dimakopoulou, Zhimei Ren, Zhengyuan Zhou
- abstract@[open-review](#): During online decision making in Multi-Armed Bandits (MAB), one needs to conduct inference on the true mean reward of each arm based on data collected so far at each step. However, since the arms are adaptively selected--thereby yielding non-iid data--conducting inference accurately is not straightforward. In particular, sample averaging, which is used in the family of UCB and Thompson sampling (TS) algorithms, does not provide a good choice as it suffers from bias and a lack of good statistical properties (e.g. asymptotic normality). Our thesis in this paper is that more sophisticated inference schemes that take into account the adaptive nature of the sequentially collected data can unlock further performance gains, even though both UCB and TS type algorithms are optimal in the worst case. In particular, we propose a variant of TS-style algorithms--which we call doubly adaptive TS--that leverages recent advances in causal inference and adaptively reweights the terms of a doubly robust estimator on the true mean reward of each arm. Through 20 synthetic domain experiments and a semi-synthetic experiment based on data from an A/B test of a web service, we demonstrate that using an adaptive inferential scheme (while still retaining the exploration efficacy of TS) provides clear benefits in online decision making: the

proposed DATS algorithm has superior empirical performance to existing baselines (UCB and TS) in terms of regret and sample complexity in identifying the best arm. In addition, we also provide a finite-time regret bound of doubly adaptive TS that matches (up to log factors) those of UCB and TS algorithms, thereby establishing that its improved practical benefits do not come at the expense of worst-case suboptimality.

[Efficient Truncated Linear Regression with Unknown Noise Variance](#)

- Constantinos Daskalakis, Patroklos Stefanou, Rui Yao, Emmanouil Zampetakis
- abstract@[open-review](#): Truncated linear regression is a classical challenge in Statistics, wherein a label, $y = w^T x + \epsilon$, and its corresponding feature vector, $x \in \mathbb{R}^k$, are only observed if the label falls in some subset $\epsilon \in \mathbb{R}$; otherwise the existence of the pair (x, y) is hidden from observation. Linear regression with truncated observations has remained a challenge, in its general form, since the early works of [Tobin'58, Amemiya '73]. When the distribution of the error is normal with known variance, recent work of [Daskalakis et al. '19] provides computationally and statistically efficient estimators of the linear model, w . In this paper, we provide the first computationally and statistically efficient estimators for truncated linear regression when the noise variance is unknown, estimating both the linear model and the variance of the noise. Our estimator is based on an efficient implementation of Projected Stochastic Gradient Descent on the negative log-likelihood of the truncated sample. Importantly, we show that the error of our estimates is asymptotically normal, and we use this to provide explicit confidence regions for our estimates.

[Breaking the Dilemma of Medical Image-to-image Translation](#)

- Lingke Kong, Chenyu Lian, Detian Huang, zhenjiang li, Yanle Hu, Qichao Zhou
- abstract@[open-review](#): Supervised Pix2Pix and unsupervised Cycle-consistency are two modes that dominate the field of medical image-to-image translation. However, neither modes are ideal. The Pix2Pix mode has excellent performance. But it requires paired and well pixel-wise aligned images, which may not always be achievable due to respiratory motion or anatomy change between times that paired images are acquired. The Cycle-consistency mode is less stringent with training data and works well on unpaired or misaligned images. But its performance may not be optimal. In order to break the dilemma of the existing modes, we propose a new unsupervised mode called RegGAN for medical image-to-image translation. It is based on the theory of "loss-correction". In RegGAN, the misaligned target images are considered as noisy labels and the generator is trained with an additional registration network to fit the misaligned noise distribution adaptively. The goal is to search for the common optimal solution to both image-to-image translation and registration tasks. We incorporated RegGAN into a few state-of-the-art image-to-image translation methods and demonstrated that RegGAN could be easily combined with these methods to improve their performances. Such as a simple CycleGAN in our mode surpasses latest NICEGAN even though using less network parameters. Based on our results, RegGAN outperformed both Pix2Pix on aligned data and Cycle-consistency on misaligned or unpaired data. RegGAN is insensitive to noises which makes it a better choice for a wide range of scenarios, especially for medical image-to-image translation tasks in which well pixel-wise aligned data are not available. Code and dataset are available at <https://github.com/Kid-Liet/Reg-GAN>.

[Temporally Abstract Partial Models](#)

- Khimya Khetarpal, Zafarali Ahmed, Gheorghe Comanici, Doina Precup
- abstract@[open-review](#): Humans and animals have the ability to reason and make predictions about different courses of action at many time scales. In reinforcement learning, option models (Sutton, Precup & Singh, 1999; Precup, 2000) provide the framework for this kind of temporally abstract prediction and reasoning. Natural intelligent agents are also able to focus their attention on courses of action that are relevant or feasible in a given situation, sometimes termed affordable actions. In this paper, we define a notion of affordances for options, and develop temporally abstract partial option models, that take into account the fact that an option might be affordable only in certain situations. We analyze the trade-offs between estimation and approximation error in planning and learning when using such models, and identify some interesting special cases. Additionally, we empirically demonstrate the ability to learn both affordances and partial option models online resulting in improved sample efficiency and planning time in the Taxi domain.

[TransMatcher: Deep Image Matching Through Transformers for Generalizable Person Re-identification](#)

- Shengcai Liao, Ling Shao
- abstract@[open-review](#): Transformers have recently gained increasing attention in computer vision. However, existing studies mostly use Transformers for feature representation learning, e.g. for image classification and dense predictions, and the generalizability of Transformers is unknown. In this work, we further investigate the possibility of applying Transformers for image matching and metric learning given pairs of images. We find that the Vision Transformer (ViT) and the vanilla Transformer with decoders are not adequate for image matching due to their lack of image-to-image attention. Thus, we further design two naive solutions, i.e. query-gallery concatenation in ViT, and query-gallery cross-attention in the vanilla Transformer. The latter improves the performance, but it is still limited. This implies that the attention mechanism in Transformers is primarily designed for global feature aggregation, which is not naturally suitable for image matching. Accordingly, we propose a new simplified decoder, which drops the full attention implementation with the softmax weighting, keeping only the query-key similarity computation. Additionally, global max pooling and a multilayer perceptron (MLP) head are applied to decode the matching result. This way, the simplified decoder is computationally more efficient, while at the same time more effective for image matching. The proposed method, called TransMatcher, achieves state-of-the-art performance in generalizable person re-identification, with up to 6.1% and 5.7% performance gains in Rank-1 and mAP, respectively, on several popular datasets. Code is available at <https://github.com/ShengcaiLiao/QACconv>.

[Multi-Objective SPIBB: Seldonian Offline Policy Improvement with Safety Constraints in Finite MDPs](#)

- harsh satija, Philip S. Thomas, Joelle Pineau, Romain Laroche
- abstract@[open-review](#): We study the problem of Safe Policy Improvement (SPI) under constraints in the offline Reinforcement Learning (RL) setting. We consider the scenario where: (i) we have a dataset collected under a known baseline policy, (ii) multiple reward signals are received from the environment inducing as many objectives to optimize. We present an SPI formulation for this RL setting that takes into account the preferences of the algorithm's user for handling the trade-offs for different reward signals while ensuring that the new policy performs at least as well as the baseline policy along each individual objective. We build on traditional SPI algorithms and propose a novel method based on Safe Policy Iteration with Baseline Bootstrapping (SPIBB, Laroche et al., 2019) that provides high probability guarantees on the performance of the agent in the true environment. We show the effectiveness of our method on a synthetic grid-world safety task as well as in a real-world critical care context to learn a policy for the administration of IV fluids and vasopressors to treat sepsis.

[Is Automated Topic Model Evaluation Broken? The Incoherence of Coherence](#)

- Alexander Hoyle, Pranav Goel, Andrew Hian-Cheong, Denis Peskov, Jordan Boyd-Graber, Philip Resnik
- abstract@[open-review](#): Topic model evaluation, like evaluation of other unsupervised methods, can be contentious. However, the field has coalesced around automated estimates of topic coherence, which rely on the frequency of word co-occurrences in a reference corpus. Contemporary neural topic models surpass classical ones according to these metrics. At the same time, topic model evaluation suffers from a validation gap: automated coherence, developed for classical models, has not been validated using human experimentation for neural models. In addition, a meta-analysis of topic modeling literature reveals a substantial standardization gap in automated topic modeling benchmarks. To address the validation gap, we compare automated coherence with the two most widely accepted human judgment tasks: topic rating and word intrusion. To address the standardization gap, we systematically evaluate a dominant classical model and two state-of-the-art neural models on two commonly used datasets. Automated evaluations declare

a winning model when corresponding human evaluations do not, calling into question the validity of fully automatic evaluations independent of human judgments.

[INDIGO: GNN-Based Inductive Knowledge Graph Completion Using Pair-Wise Encoding](#)

- Shuwen Liu, Bernardo Grau, Ian Horrocks, Egor Kostylev
- abstract@[open-review](#): The aim of knowledge graph (KG) completion is to extend an incomplete KG with missing triples. Popular approaches based on graph embeddings typically work by first representing the KG in a vector space, and then applying a predefined scoring function to the resulting vectors to complete the KG. These approaches work well in transductive settings, where predicted triples involve only constants seen during training; however, they are not applicable in inductive settings, where the KG on which the model was trained is extended with new constants or merged with other KGs. The use of Graph Neural Networks (GNNs) has recently been proposed as a way to overcome these limitations; however, existing approaches do not fully exploit the capabilities of GNNs and still rely on heuristics and ad-hoc scoring functions. In this paper, we propose a novel approach, where the KG is fully encoded into a GNN in a transparent way, and where the predicted triples can be read out directly from the last layer of the GNN without the need for additional components or scoring functions. Our experiments show that our model outperforms state-of-the-art approaches on inductive KG completion benchmarks.

[Do Input Gradients Highlight Discriminative Features?](#)

- Harshay Shah, Prateek Jain, Praneeth Netrapalli
- abstract@[open-review](#): Post-hoc gradient-based interpretability methods [Simonyan et al., 2013, Smilkov et al., 2017] that provide instance-specific explanations of model predictions are often based on assumption (A): magnitude of input gradients—gradients of logits with respect to input—noisily highlight discriminative task-relevant features. In this work, we test the validity of assumption (A) using a three-pronged approach:1. We develop an evaluation framework, DiffROAR, to test assumption (A) on four image classification benchmarks. Our results suggest that (i) input gradients of standard models (i.e., trained on original data) may grossly violate (A), whereas (ii) input gradients of adversarially robust models satisfy (A).2. We then introduce BlockMNIST, an MNIST-based semi-real dataset, that by design encodes a priori knowledge of discriminative features. Our analysis on BlockMNIST leverages this information to validate as well as characterize differences between input gradient attributions of standard and robust models.3. Finally, we theoretically prove that our empirical findings hold on a simplified version of the BlockMNIST dataset. Specifically, we prove that input gradients of standard one-hidden-layer MLPs trained on this dataset do not highlight instance-specific signal coordinates, thus grossly violating assumption (A).Our findings motivate the need to formalize and test common assumptions in interpretability in a falsifiable manner [Leavitt and Morcos, 2020]. We believe that the DiffROAR evaluation framework and BlockMNIST-based datasets can serve as sanity checks to audit instance-specific interpretability methods; code and data available at <https://github.com/harshays/inputgradients>.

[Improving Conditional Coverage via Orthogonal Quantile Regression](#)

- Shai Feldman, Stephen Bates, Yaniv Romano
- abstract@[open-review](#): We develop a method to generate prediction intervals that have a user-specified coverage level across all regions of feature-space, a property called conditional coverage. A typical approach to this task is to estimate the conditional quantiles with quantile regression---it is well-known that this leads to correct coverage in the large-sample limit, although it may not be accurate in finite samples. We find in experiments that traditional quantile regression can have poor conditional coverage. To remedy this, we modify the loss function to promote independence between the size of the intervals and the indicator of a miscoverage event. For the true conditional quantiles, these two quantities are independent (orthogonal), so the modified loss function continues to be valid. Moreover, we empirically show that the modified loss function leads to improved conditional coverage, as evaluated by several metrics. We also introduce two new metrics that check conditional coverage by looking at the strength of the dependence between the interval size and the indicator of miscoverage.

[Minimizing Polarization and Disagreement in Social Networks via Link Recommendation](#)

- Liwang Zhu, Qi Bao, Zhongzhi Zhang
- abstract@[open-review](#): Individual's opinions are fundamentally shaped and evolved by their interactions with other people, and social phenomena such as disagreement and polarization are now tightly woven into daily life. The quantification and optimization of these concepts have been the subject of much recent research behind a wealth of high-impact data mining applications. In particular, researchers have addressed the question of how such concepts can be optimized by influencing the opinion of a small number of individuals or by designing the network from scratch. Here, rather than a "design-from-scratch" approach or altering the initial opinion, we study the optimization problem of recommending \$k\$ new links to minimize the sum of polarization and disagreement in a social network with \$n\$ nodes and \$m\$ edges. We show that our objective function of this combinatorial optimization problem is not submodular, although it is monotone. We propose a simple greedy algorithm with a constant-factor approximation that solves the problem in cubic running time, and we provide theoretical analysis of the approximation guarantee for the algorithm. To overcome the computation challenge for large networks, we also provide a fast algorithm with computation complexity \$\mathcal{O}(n^2 \cdot m \cdot \log n)\$ for any \$\epsilon > 0\$, where the \$\mathcal{O}(\cdot)\$ notation suppresses the \$\text{poly}(\log n)\$ factors. Extensive experiments on real datasets demonstrate both the efficiency and effectiveness of our algorithms.

[Adversarial Attacks on Black Box Video Classifiers: Leveraging the Power of Geometric Transformations](#)

- Shasha Li, Abhishek Aich, Shitong Zhu, Salman Asif, Chengyu Song, Amit Roy-Chowdhury, Srikanth Krishnamurthy
- abstract@[open-review](#): When compared to the image classification models, black-box adversarial attacks against video classification models have been largely understudied. This could be possible because, with video, the temporal dimension poses significant additional challenges in gradient estimation. Query-efficient black-box attacks rely on effectively estimated gradients towards maximizing the probability of misclassifying the target video. In this work, we demonstrate that such effective gradients can be searched for by parameterizing the temporal structure of the search space with geometric transformations. Specifically, we design a novel iterative algorithm GEometric TRAnsformed Perturbations (GEO-TRAP), for attacking video classification models. GEO-TRAP employs standard geometric transformation operations to reduce the search space for effective gradients into searching for a small group of parameters that define these operations. This group of parameters describes the geometric progression of gradients, resulting in a reduced and structured search space. Our algorithm inherently leads to successful perturbations with surprisingly few queries. For example, adversarial examples generated from GEO-TRAP have better attack success rates with ~73.55% fewer queries compared to the state-of-the-art method for video adversarial attacks on the widely used Jester dataset. Overall, our algorithm exposes vulnerabilities of diverse video classification models and achieves new state-of-the-art results under black-box settings on two large datasets.

[Optimal Rates for Random Order Online Optimization](#)

- Uri Sherman, Tomer Koren, Yishay Mansour
- abstract@[open-review](#): We study online convex optimization in the random order model, recently proposed by Garber et al. (2020), where the loss functions may be chosen by an adversary, but are then presented to the online algorithm in a uniformly random order. Focusing on the scenario where the cumulative loss function is (strongly) convex, yet individual loss functions are smooth but might be non-convex, we give algorithms that achieve the optimal bounds and significantly outperform the results of Garber et al. (2020), completely removing the dimension dependence and improve their scaling with respect to the strong convexity parameter. Our analysis relies on novel connections between algorithmic stability and generalization for sampling

without-replacement analogous to those studied in the with-replacement i.i.d. setting, as well as on a refined average stability analysis of stochastic gradient descent.

[Discrete-Valued Neural Communication](#)

- Dianbo Liu, Alex M. Lamb, Kenji Kawaguchi, Anirudh Goyal ALIAS PARTH GOYAL, Chen Sun, Michael C. Mozer, Yoshua Bengio
- abstract@[open-review](#): Deep learning has advanced from fully connected architectures to structured models organized into components, e.g., the transformer composed of positional elements, modular architectures divided into slots, and graph neural nets made up of nodes. The nature of structured models is that communication among the components has a bottleneck, typically achieved by restricted connectivity and attention. In this work, we further tighten the bottleneck via discreteness of the representations transmitted between components. We hypothesize that this constraint serves as a useful form of inductive bias. Our hypothesis is motivated by past empirical work showing the benefits of discretization in non-structured architectures as well as our own theoretical results showing that discretization increases noise robustness and reduces the underlying dimensionality of the model. Building on an existing technique for discretization from the VQ-VAE, we consider multi-headed discretization with shared codebooks as the output of each architectural component. One motivating intuition is human language in which communication occurs through multiple discrete symbols. This form of communication is hypothesized to facilitate transmission of information between functional components of the brain by providing a common interlingua, just as it does for human-to-human communication. Our experiments show that discrete-valued neural communication (DVNC) substantially improves systematic generalization in a variety of architectures—transformers, modular architectures, and graph neural networks. We also show that the DVNC is robust to the choice of hyperparameters, making the method useful in practice.

[Skyformer: Remodel Self-Attention with Gaussian Kernel and Nyström Method](#)

- Yifan Chen, Qi Zeng, Heng Ji, Yun Yang
- abstract@[open-review](#): Transformers are expensive to train due to the quadratic time and space complexity in the self-attention mechanism. On the other hand, although kernel machines suffer from the same computation bottleneck in pairwise dot products, several approximation schemes have been successfully incorporated to considerably reduce their computational cost without sacrificing too much accuracy. In this work, we leverage the computation methods for kernel machines to alleviate the high computational cost and introduce Skyformer, which replaces the softmax structure with a Gaussian kernel to stabilize the model training and adapts the Nyström method to a non-positive semidefinite matrix to accelerate the computation. We further conduct theoretical analysis by showing that the matrix approximation error of our proposed method is small in the spectral norm. Experiments on Long Range Arena benchmark show that the proposed method is sufficient in getting comparable or even better performance than the full self-attention while requiring fewer computation resources.

[TransMIL: Transformer based Correlated Multiple Instance Learning for Whole Slide Image Classification](#)

- Zhuchen Shao, Hao Bian, Yang Chen, Yifeng Wang, Jian Zhang, Xiangyang Ji, yongbing zhang
- abstract@[open-review](#): Multiple instance learning (MIL) is a powerful tool to solve the weakly supervised classification in whole slide image (WSI) based pathology diagnosis. However, the current MIL methods are usually based on independent and identical distribution hypothesis, thus neglect the correlation among different instances. To address this problem, we proposed a new framework, called correlated MIL, and provided a proof for convergence. Based on this framework, we devised a Transformer based MIL (TransMIL), which explored both morphological and spatial information. The proposed TransMIL can effectively deal with unbalanced/balanced and binary/multiple classification with great visualization and interpretability. We conducted various experiments for three different computational pathology problems and achieved better performance and faster convergence compared with state-of-the-art methods. The test AUC for the binary tumor classification can be up to 93.09% over CAMELYON16 dataset. And the AUC over the cancer subtypes classification can be up to 96.03% and 98.82% over TCGA-NSCLC dataset and TCGA-RCC dataset, respectively. Implementation is available at: <https://github.com/szc19990412/TransMIL>.

[Multi-view Contrastive Graph Clustering](#)

- ErLin Pan, Zhao Kang
- abstract@[open-review](#): With the explosive growth of information technology, multi-view graph data have become increasingly prevalent and valuable. Most existing multi-view clustering techniques either focus on the scenario of multiple graphs or multi-view attributes. In this paper, we propose a generic framework to cluster multi-view attributed graph data. Specifically, inspired by the success of contrastive learning, we propose multi-view contrastive graph clustering (MCGC) method to learn a consensus graph since the original graph could be noisy or incomplete and is not directly applicable. Our method composes of two key steps: we first filter out the undesirable high-frequency noise while preserving the graph geometric features via graph filtering and obtain a smooth representation of nodes; we then learn a consensus graph regularized by graph contrastive loss. Results on several benchmark datasets show the superiority of our method with respect to state-of-the-art approaches. In particular, our simple approach outperforms existing deep learning-based methods.

[Inverse-Weighted Survival Games](#)

- Xintian Han, Mark Goldstein, Aahlad Puli, Thomas Wies, Adler Perotte, Rajesh Ranganath
- abstract@[open-review](#): Deep models trained through maximum likelihood have achieved state-of-the-art results for survival analysis. Despite this training scheme, practitioners evaluate models under other criteria, such as binary classification losses at a chosen set of time horizons, e.g. Brier score (BS) and Bernoulli log likelihood (BLL). Models trained with maximum likelihood may have poor BS or BLL since maximum likelihood does not directly optimize these criteria. Directly optimizing criteria like BS requires inverse-weighting by the censoring distribution. However, estimating the censoring model under these metrics requires inverse-weighting by the failure distribution. The objective for each model requires the other, but neither are known. To resolve this dilemma, we introduce Inverse-Weighted Survival Games. In these games, objectives for each model are built from re-weighted estimates featuring the other model, where the latter is held fixed during training. When the loss is proper, we show that the games always have the true failure and censoring distributions as a stationary point. This means models in the game do not leave the correct distributions once reached. We construct one case where this stationary point is unique. We show that these games optimize BS on simulations and then apply these principles on real world cancer and critically-ill patient data.

[Generalization Bounds for Meta-Learning via PAC-Bayes and Uniform Stability](#)

- Alec Farid, Anirudha Majumdar
- abstract@[open-review](#): We are motivated by the problem of providing strong generalization guarantees in the context of meta-learning. Existing generalization bounds are either challenging to evaluate or provide vacuous guarantees in even relatively simple settings. We derive a probably approximately correct (PAC) bound for gradient-based meta-learning using two different generalization frameworks in order to deal with the qualitatively different challenges of generalization at the "base" and "meta" levels. We employ bounds for uniformly stable algorithms at the base level and bounds from the PAC-Bayes framework at the meta level. The result of this approach is a novel PAC bound that is tighter when the base learner adapts quickly, which is precisely the goal of meta-learning. We show that our bound provides a tighter guarantee than other bounds on a toy non-convex problem on the unit sphere and a text-based classification example. We also present a practical regularization scheme motivated by the bound in settings where the bound is loose and demonstrate improved performance over baseline techniques.

Parallel Bayesian Optimization of Multiple Noisy Objectives with Expected Hypervolume Improvement

- Samuel Daulton, Maximilian Balandat, Eytan Bakshy
- abstract@[open-review](#): Optimizing multiple competing black-box objectives is a challenging problem in many fields, including science, engineering, and machine learning. Multi-objective Bayesian optimization (MOBO) is a sample-efficient approach for identifying the optimal trade-offs between the objectives. However, many existing methods perform poorly when the observations are corrupted by noise. We propose a novel acquisition function, NEHVI, that overcomes this important practical limitation by applying a Bayesian treatment to the popular expected hypervolume improvement (EHVI) criterion and integrating over this uncertainty in the Pareto frontier. We argue that, even in the noiseless setting, generating multiple candidates in parallel is an incarnation of EHVI with uncertainty in the Pareto frontier and therefore can be addressed using the same underlying technique. Through this lens, we derive a natural parallel variant, qNEHVI, that reduces computational complexity of parallel EHVI from exponential to polynomial with respect to the batch size. qNEHVI is one-step Bayes-optimal for hypervolume maximization in both noisy and noiseless environments, and we show that it can be optimized effectively with gradient-based methods via sample average approximation. Empirically, we demonstrate not only that qNEHVI is substantially more robust to observation noise than existing MOBO approaches, but also that it achieves state-of-the-art optimization performance and competitive wall-times in large-batch environments.

Evolution Gym: A Large-Scale Benchmark for Evolving Soft Robots

- Jagdeep Bhatia, Holly Jackson, Yunsheng Tian, Jie Xu, Wojciech Matusik
- abstract@[open-review](#): Both the design and control of a robot play equally important roles in its task performance. However, while optimal control is well studied in the machine learning and robotics community, less attention is placed on finding the optimal robot design. This is mainly because co-optimizing design and control in robotics is characterized as a challenging problem, and more importantly, a comprehensive evaluation benchmark for co-optimization does not exist. In this paper, we propose Evolution Gym, the first large-scale benchmark for co-optimizing the design and control of soft robots. In our benchmark, each robot is composed of different types of voxels (e.g., soft, rigid, actuators), resulting in a modular and expressive robot design space. Our benchmark environments span a wide range of tasks, including locomotion on various types of terrains and manipulation. Furthermore, we develop several robot co-evolution algorithms by combining state-of-the-art design optimization methods and deep reinforcement learning techniques. Evaluating the algorithms on our benchmark platform, we observe robots exhibiting increasingly complex behaviors as evolution progresses, with the best evolved designs solving many of our proposed tasks. Additionally, even though robot designs are evolved autonomously from scratch without prior knowledge, they often grow to resemble existing natural creatures while outperforming hand-designed robots. Nevertheless, all tested algorithms fail to find robots that succeed in our hardest environments. This suggests that more advanced algorithms are required to explore the high-dimensional design space and evolve increasingly intelligent robots -- an area of research in which we hope Evolution Gym will accelerate progress. Our website with code, environments, documentation, and tutorials is available at <http://evogym.csail.mit.edu/>.

On Calibration and Out-of-Domain Generalization

- Yoav Wald, Amir Feder, Daniel Greenfeld, Uri Shalit
- abstract@[open-review](#): Out-of-domain (OOD) generalization is a significant challenge for machine learning models. Many techniques have been proposed to overcome this challenge, often focused on learning models with certain invariance properties. In this work, we draw a link between OOD performance and model calibration, arguing that calibration across multiple domains can be viewed as a special case of an invariant representation leading to better OOD generalization. Specifically, we show that under certain conditions, models which achieve \emph{multi-domain calibration} are provably free of spurious correlations. This leads us to propose multi-domain calibration as a measurable and trainable surrogate for the OOD performance of a classifier. We therefore introduce methods that are easy to apply and allow practitioners to improve multi-domain calibration by training or modifying an existing model, leading to better performance on unseen domains. Using four datasets from the recently proposed WILDS OOD benchmark, as well as the Colored MNIST, we demonstrate that training or tuning models so they are calibrated across multiple domains leads to significantly improved performance on unseen test domains. We believe this intriguing connection between calibration and OOD generalization is promising from both a practical and theoretical point of view.

On the Convergence and Sample Efficiency of Variance-Reduced Policy Gradient Method

- Junyu Zhang, Chengzhuo Ni, zheng Yu, Csaba Szepesvari, Mengdi Wang
- abstract@[open-review](#): Policy gradient (PG) gives rise to a rich class of reinforcement learning (RL) methods. Recently, there has been an emerging trend to augment the existing PG methods such as REINFORCE by the \emph{variance reduction} techniques. However, all existing variance-reduced PG methods heavily rely on an uncheckable importance weight assumption made for every single iteration of the algorithms. In this paper, a simple gradient truncation mechanism is proposed to address this issue. Moreover, we design a Truncated Stochastic Incremental Variance-Reduced Policy Gradient (TSIVR-PG) method, which is able to maximize not only a cumulative sum of rewards but also a general utility function over a policy's long-term visiting distribution. We show an $\tilde{\mathcal{O}}(\epsilon^{-3})$ sample complexity for TSIVR-PG to find an ϵ -stationary policy. By assuming the \emph{overparameterization} of policy and exploiting the \emph{hidden convexity} of the problem, we further show that TSIVR-PG converges to global ϵ -optimal policy with $\tilde{\mathcal{O}}(\epsilon^{-2})$ samples.

Circa: Stochastic ReLUs for Private Deep Learning

- Zahra Ghodsi, Nandan Kumar Jha, Brandon Reagen, Siddharth Garg
- abstract@[open-review](#): The simultaneous rise of machine learning as a service and concerns over user privacy have increasingly motivated the need for private inference (PI). While recent work demonstrates PI is possible using cryptographic primitives, the computational overheads render it impractical. State-of-art deep networks are inadequate in this context because the source of slowdown in PI stems from the ReLU operations whereas optimizations for plaintext inference focus on reducing FLOPs. In this paper we re-think ReLU computations and propose optimizations for PI tailored to properties of neural networks. Specifically, we reformulate ReLU as an approximate sign test and introduce a novel truncation method for the sign test that significantly reduces the cost per ReLU. These optimizations result in a specific type of stochastic ReLU. The key observation is that the stochastic fault behavior is well suited for the fault-tolerant properties of neural network inference. Thus, we provide significant savings without impacting accuracy. We collectively call the optimizations Circa and demonstrate improvements of up to 4.7\$ times\$ storage and 3\$ times\$ runtime over baseline implementations; we further show that Circa can be used on top of recent PI optimizations to obtain 1.8\$ times\$ additional speedup.

Reinforcement Learning in Reward-Mixing MDPs

- Jeongyeol Kwon, Yonathan Efroni, Constantine Caramanis, Shie Mannor
- abstract@[open-review](#): Learning a near optimal policy in a partially observable system remains an elusive challenge in contemporary reinforcement learning. In this work, we consider episodic reinforcement learning in a reward-mixing Markov decision process (MDP). There, a reward function is drawn from one of S^H possible reward models at the beginning of every episode, but the identity of the chosen reward model is not revealed to the agent. Hence, the latent state space, for which the dynamics are Markovian, is not given to the agent. We study the problem of learning a near optimal policy for two reward-mixing MDPs. Unlike existing approaches that rely on strong assumptions on the dynamics, we make no assumptions and study the problem in full generality. Indeed, with no further assumptions, even for two switching reward-models, the problem requires several new ideas beyond existing algorithmic and analysis techniques for efficient exploration. We provide the first polynomial-time algorithm that finds an ϵ -optimal policy after exploring $\tilde{\mathcal{O}}(poly(H, \epsilon^{-1}) \cdot S^2 A^2)$ episodes, where H is time-horizon and S, A are the number of states and

actions respectively. This is the first efficient algorithm that does not require any assumptions in partially observed environments where the observation space is smaller than the latent state space.

[A Gang of Adversarial Bandits](#)

- Mark Herbster, Stephen Pasteris, Fabio Vitale, Massimiliano Pontil
- abstract@[open-review](#): We consider running multiple instances of multi-armed bandit (MAB) problems in parallel. A main motivation for this study are online recommendation systems, in which each of $\$N\$$ users is associated with a MAB problem and the goal is to exploit users' similarity in order to learn users' preferences to $\$K\$$ items more efficiently. We consider the adversarial MAB setting, whereby an adversary is free to choose which user and which loss to present to the learner during the learning process. Users are in a social network and the learner is aided by a-priori knowledge of the strengths of the social links between all pairs of users. It is assumed that if the social link between two users is strong then they tend to share the same action. The regret is measured relative to an arbitrary function which maps users to actions. The smoothness of the function is captured by a resistance-based dispersion measure $\$Psi\$$. We present two learning algorithms, GABA-I and GABA-II, which exploit the network structure to bias towards functions of low $\$Psi\$$ values. We show that GABA-I has an expected regret bound of $\$O(\sqrt{\ln(NK\Psi)\Psi KT})\$$ and per-trial time complexity of $\$O(K\ln(N))\$$, whilst GABA-II has a weaker $\$O(\sqrt{\ln(N\Psi)\ln(NK\Psi)\Psi KT})\$$ regret, but a better $\$O(\ln(K)\ln(N))\$$ per-trial time complexity. We highlight improvements of both algorithms over running independent standard MABs across users.

[Explaining Hyperparameter Optimization via Partial Dependence Plots](#)

- Julia Moosbauer, Julia Herbinger, Giuseppe Casalicchio, Marius Lindauer, Bernd Bischl
- abstract@[open-review](#): Automated hyperparameter optimization (HPO) can support practitioners to obtain peak performance in machine learning models. However, there is often a lack of valuable insights into the effects of different hyperparameters on the final model performance. This lack of explainability makes it difficult to trust and understand the automated HPO process and its results. We suggest using interpretable machine learning (IML) to gain insights from the experimental data obtained during HPO with Bayesian optimization (BO). BO tends to focus on promising regions with potential high-performance configurations and thus induces a sampling bias. Hence, many IML techniques, such as the partial dependence plot (PDP), carry the risk of generating biased interpretations. By leveraging the posterior uncertainty of the BO surrogate model, we introduce a variant of the PDP with estimated confidence bands. We propose to partition the hyperparameter space to obtain more confident and reliable PDPs in relevant sub-regions. In an experimental study, we provide quantitative evidence for the increased quality of the PDPs within sub-regions.

[Robustifying Algorithms of Learning Latent Trees with Vector Variables](#)

- Fengzhuo Zhang, Vincent Tan
- abstract@[open-review](#): We consider learning the structures of Gaussian latent tree models with vector observations when a subset of them are arbitrarily corrupted. First, we present the sample complexities of Recursive Grouping (RG) and Chow-Liu Recursive Grouping (CLRG) without the assumption that the effective depth is bounded in the number of observed nodes, significantly generalizing the results in Choi et al. (2011). We show that Chow-Liu initialization in CLRG greatly reduces the sample complexity of RG from being exponential in the diameter of the tree to only logarithmic in the diameter for the hidden Markov model (HMM). Second, we robustify RG, CLRG, Neighbor Joining (NJ) and Spectral NJ (SNJ) by using the truncated inner product. These robustified algorithms can tolerate a number of corruptions up to the square root of the number of clean samples. Finally, we derive the first known instance-dependent impossibility result for structure learning of latent trees. The optimalities of the robust version of CLRG and NJ are verified by comparing their sample complexities and the impossibility result.

[Representation Learning on Spatial Networks](#)

- Zheng Zhang, Liang Zhao
- abstract@[open-review](#): Spatial networks are networks for which the nodes and edges are constrained by geometry and embedded in real space, which has crucial effects on their topological properties. Although tremendous success has been achieved in spatial and network representation separately in recent years, there exist very little works on the representation of spatial networks. Extracting powerful representations from spatial networks requires the development of appropriate tools to uncover the pairing of both spatial and network information in the appearance of node permutation invariant, and rotation and translation invariant. Hence it can not be modeled merely with either spatial or network models individually. To address these challenges, this paper proposes a generic framework for spatial network representation learning. Specifically, a provably information-lossless and roto-translation invariant representation of spatial information on networks is presented. Then a higher-order spatial network convolution operation that adapts to our proposed representation is introduced. To ensure efficiency, we also propose a new approach that relied on sampling random spanning trees to reduce the time and memory complexity from $\$O(N^3)\$$ to $\$O(N)\$$. We demonstrate the strength of our proposed framework through extensive experiments on both synthetic and real-world datasets. The code for the proposed model is available at \url{https://github.com/rollingstonezz/SGMP_code}.

[Continuous-time edge modelling using non-parametric point processes](#)

- Xuhui Fan, Bin Li, Feng Zhou, Scott Sisson
- abstract@[open-review](#): The mutually-exciting Hawkes process (ME-HP) is a natural choice to model reciprocity, which is an important attribute of continuous-time edge (dyadic) data. However, existing ways of implementing the ME-HP for such data are either inflexible, as the exogenous (background) rate functions are typically constant and the endogenous (excitation) rate functions are specified parametrically, or inefficient, as inference usually relies on Markov chain Monte Carlo methods with high computational costs. To address these limitations, we discuss various approaches to model design, and develop three variants of non-parametric point processes for continuous-time edge modelling (CTEM). The resulting models are highly adaptable as they generate intensity functions through sigmoidal Gaussian processes, and so provide greater modelling flexibility than parametric forms. The models are implemented via a fast variational inference method enabled by a novel edge modelling construction. The superior performance of the proposed CTEM models is demonstrated through extensive experimental evaluations on four real-world continuous-time edge data sets.

[Deep inference of latent dynamics with spatio-temporal super-resolution using selective backpropagation through time](#)

- Feng Zhu, Andrew Sedler, Harrison A Grier, Nauman Ahad, Mark Davenport, Matthew Kaufman, Andrea Giovannucci, Chethan Pandarinath
- abstract@[open-review](#): Modern neural interfaces allow access to the activity of up to a million neurons within brain circuits. However, bandwidth limits often create a trade-off between greater spatial sampling (more channels or pixels) and the temporal frequency of sampling. Here we demonstrate that it is possible to obtain spatio-temporal super-resolution in neuronal time series by exploiting relationships among neurons, embedded in latent low-dimensional population dynamics. Our novel neural network training strategy, selective backpropagation through time (SBTT), enables learning of deep generative models of latent dynamics from data in which the set of observed variables changes at each time step. The resulting models are able to infer activity for missing samples by combining observations with learned latent dynamics. We test SBTT applied to sequential autoencoders and demonstrate more efficient and higher-fidelity characterization of neural population dynamics in electrophysiological and calcium imaging data. In electrophysiology, SBTT enables accurate inference of neuronal population dynamics with lower interface bandwidths, providing an avenue to significant power savings for implanted neuroelectronic interfaces. In applications to two-photon calcium imaging, SBTT accurately uncovers high-frequency temporal structure underlying neural population activity, substantially outperforming the current state-of-the-art. Finally, we demonstrate that performance could be further improved by using limited, high-bandwidth sampling to pretrain dynamics models, and then using SBTT to adapt these models for sparsely-sampled data.

Memory-efficient Patch-based Inference for Tiny Deep Learning

- Ji Lin, Wei-Ming Chen, Han Cai, Chuang Gan, Song Han
- abstract@[open-review](#): Tiny deep learning on microcontroller units (MCUs) is challenging due to the limited memory size. We find that the memory bottleneck is due to the imbalanced memory distribution in convolutional neural network (CNN) designs: the first several blocks have an order of magnitude larger memory usage than the rest of the network. To alleviate this issue, we propose a generic patch-by-patch inference scheduling, which operates only on a small spatial region of the feature map and significantly cuts down the peak memory. However, naive implementation brings overlapping patches and computation overhead. We further propose receptive field redistribution to shift the receptive field and FLOPs to the later stage and reduce the computation overhead. Manually redistributing the receptive field is difficult. We automate the process with neural architecture search to jointly optimize the neural architecture and inference scheduling, leading to MCUNetV2. Patch-based inference effectively reduces the peak memory usage of existing networks by 4-8x. Co-designed with neural networks, MCUNetV2 sets a record ImageNet accuracy on MCU (71.8%) and achieves >90% accuracy on the visual wake words dataset under only 32kB SRAM. MCUNetV2 also unblocks object detection on tiny devices, achieving 16.9% higher mAP on Pascal VOC compared to the state-of-the-art result. Our study largely addressed the memory bottleneck in tinyML and paved the way for various vision applications beyond image classification.

Self-Interpretable Model with Transformation Equivariant Interpretation

- Yipei Wang, Xiaoqian Wang
- abstract@[open-review](#): With the proliferation of machine learning applications in the real world, the demand for explaining machine learning predictions continues to grow especially in high-stakes fields. Recent studies have found that interpretation methods can be sensitive and unreliable, where the interpretations can be disturbed by perturbations or transformations of input data. To address this issue, we propose to learn robust interpretation through transformation equivariant regularization in a self-interpretable model. The resulting model is capable of capturing valid interpretation that is equivariant to geometric transformations. Moreover, since our model is self-interpretable, it enables faithful interpretations that reflect the true predictive mechanism. Unlike existing self-interpretable models, which usually sacrifice expressive power for the sake of interpretation quality, our model preserves the high expressive capability comparable to the state-of-the-art deep learning models in complex tasks, while providing visualizable and faithful high-quality interpretation. We compare with various related methods and validate the interpretation quality and consistency of our model.

Solving Min-Max Optimization with Hidden Structure via Gradient Descent Ascent

- Emmanouil-Vasileios Vlatakis-Gkaragkounis, Lampros Flokas, Georgios Piliouras
- abstract@[open-review](#): Many recent AI architectures are inspired by zero-sum games, however, the behavior of their dynamics is still not well understood. Inspired by this, we study standard gradient descent ascent (GDA) dynamics in a specific class of non-convex non-concave zero-sum games, that we call hidden zero-sum games. In this class, players control the inputs of smooth but possibly non-linear functions whose outputs are being applied as inputs to a convex-concave game. Unlike general zero-sum games, these games have a well-defined notion of solution; outcomes that implement the von-Neumann equilibrium of the ``hidden'' convex-concave game. We provide conditions under which vanilla GDA provably converges not merely to local Nash, but the actual von-Neumann solution. If the hidden game lacks strict convexity properties, GDA may fail to converge to any equilibrium, however, by applying standard regularization techniques we can prove convergence to a von-Neumann solution of a slightly perturbed zero-sum game. Our convergence results are non-local despite working in the setting of non-convex non-concave games. Critically, under proper assumptions we combine the Center-Stable Manifold Theorem along with novel type of initialization dependent Lyapunov functions to prove that almost all initial conditions converge to the solution. Finally, we discuss diverse applications of our framework ranging from generative adversarial networks to evolutionary biology.

Preserved central model for faster bidirectional compression in distributed settings

- Constantin Philippenko, Aymeric Dieuleveut
- abstract@[open-review](#): We develop a new approach to tackle communication constraints in a distributed learning problem with a central server. We propose and analyze a new algorithm that performs bidirectional compression and achieves the same convergence rate as algorithms using only uplink (from the local workers to the central server) compression. To obtain this improvement, we design MCM, an algorithm such that the downlink compression only impacts local models, while the global model is preserved. As a result, and contrary to previous works, the gradients on local servers are computed on perturbed models. Consequently, convergence proofs are more challenging and require a precise control of this perturbation. To ensure it, MCM additionally combines model compression with a memory mechanism. This analysis opens new doors, e.g. incorporating worker dependent randomized-models and partial participation.

Understanding Instance-based Interpretability of Variational Auto-Encoders

- Zhifeng Kong, Kamalika Chaudhuri
- abstract@[open-review](#): Instance-based interpretation methods have been widely studied for supervised learning methods as they help explain how black box neural networks predict. However, instance-based interpretations remain ill-understood in the context of unsupervised learning. In this paper, we investigate influence functions [Koh and Liang, 2017], a popular instance-based interpretation method, for a class of deep generative models called variational auto-encoders (VAE). We formally frame the counter-factual question answered by influence functions in this setting, and through theoretical analysis, examine what they reveal about the impact of training samples on classical unsupervised learning methods. We then introduce VAE-TracIn, a computationally efficient and theoretically sound solution based on Pruthi et al. [2020], for VAEs. Finally, we evaluate VAE-TracIn on several real world datasets with extensive quantitative and qualitative analysis.

Voxel-based 3D Detection and Reconstruction of Multiple Objects from a Single Image

- Feng Liu, Xiaoming Liu
- abstract@[open-review](#): Inferring 3D locations and shapes of multiple objects from a single 2D image is a long-standing objective of computer vision. Most of the existing works either predict one of these 3D properties or focus on solving both for a single object. One fundamental challenge lies in how to learn an effective representation of the image that is well-suited for 3D detection and reconstruction. In this work, we propose to learn a regular grid of 3D voxel features from the input image which is aligned with 3D scene space via a 3D feature lifting operator. Based on the 3D voxel features, our novel CenterNet-3D detection head formulates the 3D detection as keypoint detection in the 3D space. Moreover, we devise an efficient coarse-to-fine reconstruction module, including coarse-level voxelization and a novel local PCA-SDF shape representation, which enables fine detail reconstruction and two orders of magnitude faster inference than prior methods. With complementary supervision from both 3D detection and reconstruction, one enables the 3D voxel features to be geometry and context preserving, benefiting both tasks. The effectiveness of our approach is demonstrated through 3D detection and reconstruction on single-object and multiple-object scenarios.

Test-Time Classifier Adjustment Module for Model-Agnostic Domain Generalization

- Yusuke Iwasawa, Yutaka Matsuo
- abstract@[open-review](#): This paper presents a new algorithm for domain generalization (DG), \textit{test-time template adjuster} (T3A), aiming to robustify a model to unknown distribution shift. Unlike existing methods that focus on \textit{training phase}, our method focuses \textit{test phase}, i.e.,

correcting its prediction by itself during test time. Specifically, T3A adjusts a trained linear classifier (the last layer of deep neural networks) with the following procedure: (1) compute a pseudo-prototype representation for each class using online unlabeled data augmented by the base classifier trained in the source domains, (2) and then classify each sample based on its distance to the pseudo-prototypes. T3A is back-propagation-free and modifies only the linear layer; therefore, the increase in computational cost during inference is negligible and avoids the catastrophic failure might caused by stochastic optimization. Despite its simplicity, T3A can leverage knowledge about the target domain by using off-the-shelf test-time data and improve performance. We tested our method on four domain generalization benchmarks, namely PACS, VLCS, OfficeHome, and TerraIncognita, along with various backbone networks including ResNet18, ResNet50, Big Transfer (BiT), Vision Transformers (ViT), and MLP-Mixer. The results show T3A stably improves performance on unseen domains across choices of backbone networks, and outperforms existing domain generalization methods.

[Luna: Linear Unified Nested Attention](#)

- Xuezhe Ma, Xiang Kong, Sinong Wang, Chunting Zhou, Jonathan May, Hao Ma, Luke Zettlemoyer
- abstract@[open-review](#): The quadratic computational and memory complexities of the Transformer's attention mechanism have limited its scalability for modeling long sequences. In this paper, we propose Luna, a linear unified nested attention mechanism that approximates softmax attention with two nested linear attention functions, yielding only linear (as opposed to quadratic) time and space complexity. Specifically, with the first attention function, Luna packs the input sequence into a sequence of fixed length. Then, the packed sequence is unpacked using the second attention function. As compared to a more traditional attention mechanism, Luna introduces an additional sequence with a fixed length as input and an additional corresponding output, which allows Luna to perform attention operation linearly, while also storing adequate contextual information. We perform extensive evaluations on three benchmarks of sequence modeling tasks: long-context sequence modelling, neural machine translation and masked language modeling for large-scale pretraining. Competitive or even better experimental results demonstrate both the effectiveness and efficiency of Luna compared to a variety of strong baseline methods including the full-rank attention and other efficient sparse and dense attention methods.

[Iterative Causal Discovery in the Possible Presence of Latent Confounders and Selection Bias](#)

- Raanan Y. Rohekar, Shami Nisimov, Yaniv Gurwicz, Gal Novik
- abstract@[open-review](#): We present a sound and complete algorithm, called iterative causal discovery (ICD), for recovering causal graphs in the presence of latent confounders and selection bias. ICD relies on the causal Markov and faithfulness assumptions and recovers the equivalence class of the underlying causal graph. It starts with a complete graph, and consists of a single iterative stage that gradually refines this graph by identifying conditional independence (CI) between connected nodes. Independence and causal relations entailed after any iteration are correct, rendering ICD anytime. Essentially, we tie the size of the CI conditioning set to its distance on the graph from the tested nodes, and increase this value in the successive iteration. Thus, each iteration refines a graph that was recovered by previous iterations having smaller conditioning sets---a higher statistical power---which contributes to stability. We demonstrate empirically that ICD requires significantly fewer CI tests and learns more accurate causal graphs compared to FCI, FCI+, and RFCI algorithms.

[Hindsight Task Relabelling: Experience Replay for Sparse Reward Meta-RL](#)

- Charles Packer, Pieter Abbeel, Joseph E. Gonzalez
- abstract@[open-review](#): Meta-reinforcement learning (meta-RL) has proven to be a successful framework for leveraging experience from prior tasks to rapidly learn new related tasks, however, current meta-RL approaches struggle to learn in sparse reward environments. Although existing meta-RL algorithms can learn strategies for adapting to new sparse reward tasks, the actual adaptation strategies are learned using hand-shaped reward functions, or require simple environments where random exploration is sufficient to encounter sparse reward. In this paper we present a formulation of hindsight relabelling for meta-RL, which relabels experience during meta-training to enable learning to learn entirely using sparse reward. We demonstrate the effectiveness of our approach on a suite of challenging sparse reward environments that previously required dense reward during meta-training to solve. Our approach solves these environments using the true sparse reward function, with performance comparable to training with a proxy dense reward function.

[A Bayesian-Symbolic Approach to Reasoning and Learning in Intuitive Physics](#)

- Kai Xu, Akash Srivastava, Dan Gutfreund, Felix Sosa, Tomer Ullman, Josh Tenenbaum, Charles Sutton
- abstract@[open-review](#): Humans can reason about intuitive physics in fully or partially observed environments even after being exposed to a very limited set of observations. This sample-efficient intuitive physical reasoning is considered a core domain of human common sense knowledge. One hypothesis to explain this remarkable capacity, posits that humans quickly learn approximations to the laws of physics that govern the dynamics of the environment. In this paper, we propose a Bayesian-symbolic framework (BSP) for physical reasoning and learning that is close to human-level sample-efficiency and accuracy. In BSP, the environment is represented by a top-down generative model of entities, which are assumed to interact with each other under unknown force laws over their latent and observed properties. BSP models each of these entities as random variables, and uses Bayesian inference to estimate their unknown properties. For learning the unknown forces, BSP leverages symbolic regression on a novel grammar of Newtonian physics in a bilevel optimization setup. These inference and regression steps are performed in an iterative manner using expectation-maximization, allowing BSP to simultaneously learn force laws while maintaining uncertainty over entity properties. We show that BSP is more sample-efficient compared to neural alternatives on controlled synthetic datasets, demonstrate BSP's applicability to real-world common sense scenes and study BSP's performance on tasks previously used to study human physical reasoning.

[Associating Objects with Transformers for Video Object Segmentation](#)

- Zongxin Yang, Yunchao Wei, Yi Yang
- abstract@[open-review](#): This paper investigates how to realize better and more efficient embedding learning to tackle the semi-supervised video object segmentation under challenging multi-object scenarios. The state-of-the-art methods learn to decode features with a single positive object and thus have to match and segment each target separately under multi-object scenarios, consuming multiple times computing resources. To solve the problem, we propose an Associating Objects with Transformers (AOT) approach to match and decode multiple objects uniformly. In detail, AOT employs an identification mechanism to associate multiple targets into the same high-dimensional embedding space. Thus, we can simultaneously process multiple objects' matching and segmentation decoding as efficiently as processing a single object. For sufficiently modeling multi-object association, a Long Short-Term Transformer is designed for constructing hierarchical matching and propagation. We conduct extensive experiments on both multi-object and single-object benchmarks to examine AOT variant networks with different complexities. Particularly, our R50-AOT-L outperforms all the state-of-the-art competitors on three popular benchmarks, i.e., YouTube-VOS (84.1% J&F), DAVIS 2017 (84.9%), and DAVIS 2016 (91.1%), while keeping more than 3X faster multi-object run-time. Meanwhile, our AOT-T can maintain real-time multi-object speed on the above benchmarks. Based on AOT, we ranked 1st in the 3rd Large-scale VOS Challenge.

[Automatic Symmetry Discovery with Lie Algebra Convolutional Network](#)

- Nima Dehmamy, Robin Walters, Yanchen Liu, Dashun Wang, Rose Yu
- abstract@[open-review](#): Existing equivariant neural networks require prior knowledge of the symmetry group and discretization for continuous groups. We propose to work with Lie algebras (infinitesimal generators) instead of Lie groups. Our model, the Lie algebra convolutional network (L-conv) can automatically discover symmetries and does not require discretization of the group. We show that L-conv can serve as a building block to construct any

group equivariant feedforward architecture. Both CNNs and Graph Convolutional Networks can be expressed as L-conv with appropriate groups. We discover direct connections between L-conv and physics: (1) group invariant loss generalizes field theory (2) Euler-Lagrange equation measures the robustness, and (3) equivariance leads to conservation laws and Noether current. These connections open up new avenues for designing more general equivariant networks and applying them to important problems in physical sciences.

[Zero Time Waste: Recycling Predictions in Early Exit Neural Networks](#)

- Maciej Wołczyk, Bartosz Wójcik, Klaudia Bałazy, Igor T Podolak, Jacek Tabor, Marek Śmieja, Tomasz Trzcinski
- abstract@[open-review](#): The problem of reducing processing time of large deep learning models is a fundamental challenge in many real-world applications. Early exit methods strive towards this goal by attaching additional Internal Classifiers (ICs) to intermediate layers of a neural network. ICs can quickly return predictions for easy examples and, as a result, reduce the average inference time of the whole model. However, if a particular IC does not decide to return an answer early, its predictions are discarded, with its computations effectively being wasted. To solve this issue, we introduce Zero Time Waste (ZTW), a novel approach in which each IC reuses predictions returned by its predecessors by (1) adding direct connections between ICs and (2) combining previous outputs in an ensemble-like manner. We conduct extensive experiments across various datasets and architectures to demonstrate that ZTW achieves a significantly better accuracy vs. inference time trade-off than other recently proposed early exit methods.

[On Model Calibration for Long-Tailed Object Detection and Instance Segmentation](#)

- Tai-Yu Pan, Cheng Zhang, Yandong Li, Hexiang Hu, Dong Xuan, Soravit Changpinyo, Boqing Gong, Wei-Lun Chao
- abstract@[open-review](#): Vanilla models for object detection and instance segmentation suffer from the heavy bias toward detecting frequent objects in the long-tailed setting. Existing methods address this issue mostly during training, e.g., by re-sampling or re-weighting. In this paper, we investigate a largely overlooked approach --- post-processing calibration of confidence scores. We propose NorCal, Normalized Calibration for long-tailed object detection and instance segmentation, a simple and straightforward recipe that reweights the predicted scores of each class by its training sample size. We show that separately handling the background class and normalizing the scores over classes for each proposal are keys to achieving superior performance. On the LVIS dataset, NorCal can effectively improve nearly all the baseline models not only on rare classes but also on common and frequent classes. Finally, we conduct extensive analysis and ablation studies to offer insights into various modeling choices and mechanisms of our approach. Our code is publicly available at <https://github.com/tydpan/NorCal>.

[ReSSL: Relational Self-Supervised Learning with Weak Augmentation](#)

- Mingkai Zheng, Shan You, Fei Wang, Chen Qian, Changshui Zhang, Xiaogang Wang, Chang Xu
- abstract@[open-review](#): Self-supervised Learning (SSL) including the mainstream contrastive learning has achieved great success in learning visual representations without data annotations. However, most of methods mainly focus on the instance level information (ie, the different augmented images of the same instance should have the same feature or cluster into the same class), but there is a lack of attention on the relationships between different instances. In this paper, we introduced a novel SSL paradigm, which we term as relational self-supervised learning (ReSSL) framework that learns representations by modeling the relationship between different instances. Specifically, our proposed method employs sharpened distribution of pairwise similarities among different instances as $\text{textit\{relation\}}$ metric, which is thus utilized to match the feature embeddings of different augmentations. Moreover, to boost the performance, we argue that weak augmentations matter to represent a more reliable relation, and leverage momentum strategy for practical efficiency. Experimental results show that our proposed ReSSL significantly outperforms the previous state-of-the-art algorithms in terms of both performance and training efficiency.

[Learning to See by Looking at Noise](#)

- Manel Baradad Jurjo, Jonas Wulff, Tongzhou Wang, Phillip Isola, Antonio Torralba
- abstract@[open-review](#): Current vision systems are trained on huge datasets, and these datasets come with costs: curation is expensive, they inherit human biases, and there are concerns over privacy and usage rights. To counter these costs, interest has surged in learning from cheaper data sources, such as unlabeled images. In this paper we go a step further and ask if we can do away with real image datasets entirely, instead learning from procedural noise processes. We investigate a suite of image generation models that produce images from simple random processes. These are then used as training data for a visual representation learner with a contrastive loss. In particular, we study statistical image models, randomly initialized deep generative models, and procedural graphics models. Our findings show that it is important for the noise to capture certain structural properties of real data but that good performance can be achieved even with processes that are far from realistic. We also find that diversity is a key property to learn good representations.

[Explicit loss asymptotics in the gradient descent training of neural networks](#)

- Maksim Velikanov, Dmitry Yarotsky
- abstract@[open-review](#): Current theoretical results on optimization trajectories of neural networks trained by gradient descent typically have the form of rigorous but potentially loose bounds on the loss values. In the present work we take a different approach and show that the learning trajectory of a wide network in a lazy training regime can be characterized by an explicit asymptotic at large training times. Specifically, the leading term in the asymptotic expansion of the loss behaves as a power law $L(t) \sim C t^{-\xi}$ with exponent ξ expressed only through the data dimension, the smoothness of the activation function, and the class of function being approximated. Our results are based on spectral analysis of the integral operator representing the linearized evolution of a large network trained on the expected loss. Importantly, the techniques we employ do not require a specific form of the data distribution, for example Gaussian, thus making our findings sufficiently universal.

[Test-Time Personalization with a Transformer for Human Pose Estimation](#)

- Yizhuo Li, Miao Hao, Zonglin Di, Nitesh Bharadwaj Gundavarapu, Xiaolong Wang
- abstract@[open-review](#): We propose to personalize a 2D human pose estimator given a set of test images of a person without using any manual annotations. While there is a significant advancement in human pose estimation, it is still very challenging for a model to generalize to different unknown environments and unseen persons. Instead of using a fixed model for every test case, we adapt our pose estimator during test time to exploit person-specific information. We first train our model on diverse data with both a supervised and a self-supervised pose estimation objectives jointly. We use a Transformer model to build a transformation between the self-supervised keypoints and the supervised keypoints. During test time, we personalize and adapt our model by fine-tuning with the self-supervised objective. The pose is then improved by transforming the updated self-supervised keypoints. We experiment with multiple datasets and show significant improvements on pose estimations with our self-supervised personalization. Project page with code is available at <https://liyz15.github.io/TTP/>.

[Towards Scalable Unpaired Virtual Try-On via Patch-Routed Spatially-Adaptive GAN](#)

- Zhenyu Xie, Zaiyu Huang, Fuwei Zhao, Haoye Dong, Michael Kampffmeyer, Xiaodan Liang
- abstract@[open-review](#): Image-based virtual try-on is one of the most promising applications of human-centric image generation due to its tremendous real-world potential. Yet, as most try-on approaches fit in-shop garments onto a target person, they require the laborious and restrictive construction of a paired training dataset, severely limiting their scalability. While a few recent works attempt to transfer garments directly from one person to another,

alleviating the need to collect paired datasets, their performance is impacted by the lack of paired (supervised) information. In particular, disentangling style and spatial information of the garment becomes a challenge, which existing methods either address by requiring auxiliary data or extensive online optimization procedures, thereby still inhibiting their scalability. To achieve a scalable virtual try-on system that can transfer arbitrary garments between a source and a target person in an unsupervised manner, we thus propose a texture-preserving end-to-end network, the PAatch-routed SpaTially-Adaptive GAN (PASTA-GAN), that facilitates real-world unpaired virtual try-on. Specifically, to disentangle the style and spatial information of each garment, PASTA-GAN consists of an innovative patch-routed disentanglement module for successfully retaining garment texture and shape characteristics. Guided by the source person's keypoints, the patch-routed disentanglement module first decouples garments into normalized patches, thus eliminating the inherent spatial information of the garment, and then reconstructs the normalized patches to the warped garment complying with the target person pose. Given the warped garment, PASTA-GAN further introduces novel spatially-adaptive residual blocks that guide the generator to synthesize more realistic garment details. Extensive comparisons with paired and unpaired approaches demonstrate the superiority of PASTA-GAN, highlighting its ability to generate high-quality try-on images when faced with a large variety of garments(e.g. vests, shirts, pants), taking a crucial step towards real-world scalable try-on.

[Bias Out-of-the-Box: An Empirical Analysis of Intersectional Occupational Biases in Popular Generative Language Models](#)

- Hannah Rose Kirk, Yennie Jun, Filippo Volpin, Haider Iqbal, Elias Benussi, Frederic Dreyer, Aleksandar Shtedritski, Yuki Asano
- abstract@[open-review](#): The capabilities of natural language models trained on large-scale data have increased immensely over the past few years. Open source libraries such as HuggingFace have made these models easily available and accessible. While prior research has identified biases in large language models, this paper considers biases contained in the most popular versions of these models when applied `out-of-the-box' for downstream tasks. We focus on generative language models as they are well-suited for extracting biases inherited from training data. Specifically, we conduct an in-depth analysis of GPT-2, which is the most downloaded text generation model on HuggingFace, with over half a million downloads per month. We assess biases related to occupational associations for different protected categories by intersecting gender with religion, sexuality, ethnicity, political affiliation, and continental name origin. Using a template-based data collection pipeline, we collect 396K sentence completions made by GPT-2 and find: (i) The machine-predicted jobs are less diverse and more stereotypical for women than for men, especially for intersections; (ii) Intersectional interactions are highly relevant for occupational associations, which we quantify by fitting 262 logistic models; (iii) For most occupations, GPT-2 reflects the skewed gender and ethnicity distribution found in US Labor Bureau data, and even pulls the societally-skewed distribution towards gender parity in cases where its predictions deviate from real labor market observations. This raises the normative question of what language models \textit{should} learn - whether they should reflect or correct for existing inequalities.

[Weisfeiler and Lehman Go Cellular: CW Networks](#)

- Cristian Bodnar, Fabrizio Frasca, Nina Otter, Yuguang Wang, Pietro Liò, Guido F. Montufar, Michael Bronstein
- abstract@[open-review](#): Graph Neural Networks (GNNs) are limited in their expressive power, struggle with long-range interactions and lack a principled way to model higher-order structures. These problems can be attributed to the strong coupling between the computational graph and the input graph structure. The recently proposed Message Passing Simplicial Networks naturally decouple these elements by performing message passing on the clique complex of the graph. Nevertheless, these models can be severely constrained by the rigid combinatorial structure of Simplicial Complexes (SCs). In this work, we extend recent theoretical results on SCs to regular Cell Complexes, topological objects that flexibly subsume SCs and graphs. We show that this generalisation provides a powerful set of graph "lifting" transformations, each leading to a unique hierarchical message passing procedure. The resulting methods, which we collectively call CW Networks (CWNs), are strictly more powerful than the WL test and not less powerful than the 3-WL test. In particular, we demonstrate the effectiveness of one such scheme, based on rings, when applied to molecular graph problems. The proposed architecture benefits from provably larger expressivity than commonly used GNNs, principled modelling of higher-order signals and from compressing the distances between nodes. We demonstrate that our model achieves state-of-the-art results on a variety of molecular datasets.

[Learning Conjoint Attentions for Graph Neural Nets](#)

- Tiantian He, Yew Soon Ong, L Bai
- abstract@[open-review](#): In this paper, we present Conjoint Attentions (CAs), a class of novel learning-to-attend strategies for graph neural networks (GNNs). Besides considering the layer-wise node features propagated within the GNN, CAs can additionally incorporate various structural interventions, such as node cluster embedding, and higher-order structural correlations that can be learned outside of GNN, when computing attention scores. The node features that are regarded as significant by the conjoint criteria are therefore more likely to be propagated in the GNN. Given the novel Conjoint Attention strategies, we then propose Graph conjoint attention networks (CATs) that can learn representations embedded with significant latent features deemed by the Conjoint Attentions. Besides, we theoretically validate the discriminative capacity of CATs. CATs utilizing the proposed Conjoint Attention strategies have been extensively tested in well-established benchmarking datasets and comprehensively compared with state-of-the-art baselines. The obtained notable performance demonstrates the effectiveness of the proposed Conjoint Attentions.

[Hybrid Regret Bounds for Combinatorial Semi-Bandits and Adversarial Linear Bandits](#)

- Shinji Ito
- abstract@[open-review](#): This study aims to develop bandit algorithms that automatically exploit tendencies of certain environments to improve performance, without any prior knowledge regarding the environments. We first propose an algorithm for combinatorial semi-bandits with a hybrid regret bound that includes two main features: a best-of-three-worlds guarantee and multiple data-dependent regret bounds. The former means that the algorithm will work nearly optimally in all environments in an adversarial setting, a stochastic setting, or a stochastic setting with adversarial corruptions. The latter implies that, even if the environment is far from exhibiting stochastic behavior, the algorithm will perform better as long as the environment is "easy" in terms of certain metrics. The metrics w.r.t. the easiness referred to in this paper include cumulative loss for optimal actions, total quadratic variation of losses, and path-length of a loss sequence. We also show hybrid data-dependent regret bounds for adversarial linear bandits, which include a first path-length regret bound that is tight up to logarithmic factors.

[Pay Better Attention to Attention: Head Selection in Multilingual and Multi-Domain Sequence Modeling](#)

- Hongyu Gong, Yun Tang, Juan Pino, Xian Li
- abstract@[open-review](#): Multi-head attention has each of the attention heads collect salient information from different parts of an input sequence, making it a powerful mechanism for sequence modeling. Multilingual and multi-domain learning are common scenarios for sequence modeling, where the key challenge is to maximize positive transfer and mitigate negative interference across languages and domains. In this paper, we find that non-selective attention sharing is sub-optimal for achieving good generalization across all languages and domains. We further propose attention sharing strategies to facilitate parameter sharing and specialization in multilingual and multi-domain sequence modeling. Our approach automatically learns shared and specialized attention heads for different languages and domains. Evaluated in various tasks including speech recognition, text-to-text and speech-to-text translation, the proposed attention sharing strategies consistently bring gains to sequence models built upon multi-head attention. For speech-to-text translation, our approach yields an average of \$+2.0\$ BLEU over \$13\$ language directions in multilingual setting and \$+2.0\$ BLEU over \$3\$ domains in multi-domain setting.

[Cardinality-Regularized Hawkes-Granger Model](#)

- Tsuyoshi Ide, Georgios Kollias, Dzung Phan, Naoki Abe

- abstract@[open-review](#): We propose a new sparse Granger-causal learning framework for temporal event data. We focus on a specific class of point processes called the Hawkes process. We begin by pointing out that most of the existing sparse causal learning algorithms for the Hawkes process suffer from a singularity in maximum likelihood estimation. As a result, their sparse solutions can appear only as numerical artifacts. In this paper, we propose a mathematically well-defined sparse causal learning framework based on a cardinality-regularized Hawkes process, which remedies the pathological issues of existing approaches. We leverage the proposed algorithm for the task of instance-wise causal event analysis, where sparsity plays a critical role. We validate the proposed framework with two real use-cases, one from the power grid and the other from the cloud data center management domain.

[Aligned Structured Sparsity Learning for Efficient Image Super-Resolution](#)

- Yulun Zhang, Huan Wang, Can Qin, Yun Fu
- abstract@[open-review](#): Lightweight image super-resolution (SR) networks have obtained promising results with moderate model size. Many SR methods have focused on designing lightweight architectures, which neglect to further reduce the redundancy of network parameters. On the other hand, model compression techniques, like neural architecture search and knowledge distillation, typically consume considerable memory and computation resources. In contrast, network pruning is a cheap and effective model compression technique. However, it is hard to be applied to SR networks directly, because filter pruning for residual blocks is well-known tricky. To address the above issues, we propose aligned structured sparsity learning (ASSL), which introduces a weight normalization layer and applies $\$L_2\$$ regularization to the scale parameters for sparsity. To align the pruned locations across different layers, we propose a \{sparsity structure alignment\} penalty term, which minimizes the norm of soft mask gram matrix. We apply aligned structured sparsity learning strategy to train efficient image SR network, named as ASSLN, with smaller model size and lower computation than state-of-the-art methods. We conduct extensive comparisons with lightweight SR networks. Our ASSLN achieves superior performance gains over recent methods quantitatively and visually.

[Why Lottery Ticket Wins? A Theoretical Perspective of Sample Complexity on Sparse Neural Networks](#)

- Shuai Zhang, Meng Wang, Sijia Liu, Pin-Yu Chen, Jinjun Xiong
- abstract@[open-review](#): The lottery ticket hypothesis (LTH) states that learning on a properly pruned network (the winning ticket) has improved test accuracy over the original unpruned network. Although LTH has been justified empirically in a broad range of deep neural network (DNN) involved applications like computer vision and natural language processing, the theoretical validation of the improved generalization of a winning ticket remains elusive. To the best of our knowledge, our work, for the first time, characterizes the performance of training a pruned neural network by analyzing the geometric structure of the objective function and the sample complexity to achieve zero generalization error. We show that the convex region near a desirable model with guaranteed generalization enlarges as the neural network model is pruned, indicating the structural importance of a winning ticket. Moreover, as the algorithm for training a pruned neural network is specified as an (accelerated) stochastic gradient descent algorithm, we theoretically show that the number of samples required for achieving zero generalization error is proportional to the number of the non-pruned weights in the hidden layer. With a fixed number of samples, training a pruned neural network enjoys a faster convergence rate to the desired model than training the original unpruned one, providing a formal justification of the improved generalization of the winning ticket. Our theoretical results are acquired from learning a pruned neural network of one hidden layer, while experimental results are further provided to justify the implications in pruning multi-layer neural networks.

[Constrained Robust Submodular Partitioning](#)

- Shengjie Wang, Tianyi Zhou, Chandrashekhar Lavania, Jeff A Bilmes
- abstract@[open-review](#): In the robust submodular partitioning problem, we aim to allocate a set of items into $\$m\$$ blocks, so that the evaluation of the minimum block according to a submodular function is maximized. Robust submodular partitioning promotes the diversity of every block in the partition. It has many applications in machine learning, e.g., partitioning data for distributed training so that the gradients computed on every block are consistent. We study an extension of the robust submodular partition problem with additional constraints (e.g., cardinality, multiple matroids, and/or knapsack) on every block. For example, when partitioning data for distributed training, we can add a constraint that the number of samples of each class is the same in each partition block, ensuring data balance. We present two classes of algorithms, i.e., Min-Block Greedy based algorithms (with an $\$O(\Omega(1/m))\$$ bound), and Round-Robin Greedy based algorithms (with a constant bound) and show that under various constraints, they still have good approximation guarantees. Interestingly, while normally the latter runs in only weakly polynomial time, we show that using the two together yields strongly polynomial running time while preserving the approximation guarantee. Lastly, we apply the algorithms on a real-world machine learning data partitioning problem showing good results.

[Online Knapsack with Frequency Predictions](#)

- Sungjin Im, Ravi Kumar, Mahshid Montazer Qaem, Manish Purohit
- abstract@[open-review](#): There has been recent interest in using machine-learned predictions to improve the worst-case guarantees of online algorithms. In this paper we continue this line of work by studying the online knapsack problem, but with very weak predictions: in the form of knowing an upper and lower bound for the number of items of each value. We systematically derive online algorithms that attain the best possible competitive ratio for any fixed prediction; we also extend the results to more general settings such as generalized one-way trading and two-stage online knapsack. Our work shows that even seemingly weak predictions can be utilized effectively to provably improve the performance of online algorithms.

[On Component Interactions in Two-Stage Recommender Systems](#)

- Jiri Hron, Karl Krauth, Michael Jordan, Niki Kilbertus
- abstract@[open-review](#): Thanks to their scalability, two-stage recommenders are used by many of today's largest online platforms, including YouTube, LinkedIn, and Pinterest. These systems produce recommendations in two steps: (i) multiple nominators—tuned for low prediction latency—preselect a small subset of candidates from the whole item pool; (ii) a slower but more accurate ranker further narrows down the nominated items, and serves to the user. Despite their popularity, the literature on two-stage recommenders is relatively scarce, and the algorithms are often treated as mere sums of their parts. Such treatment presupposes that the two-stage performance is explained by the behavior of the individual components in isolation. This is not the case: using synthetic and real-world data, we demonstrate that interactions between the ranker and the nominators substantially affect the overall performance. Motivated by these findings, we derive a generalization lower bound which shows that independent nominator training can lead to performance on par with uniformly random recommendations. We find that careful design of item pools, each assigned to a different nominator, alleviates these issues. As manual search for a good pool allocation is difficult, we propose to learn one instead using a Mixture-of-Experts based approach. This significantly improves both precision and recall at $\$K\$$.

[Lip to Speech Synthesis with Visual Context Attentional GAN](#)

- Minsu Kim, Joanna Hong, Yong Man Ro
- abstract@[open-review](#): In this paper, we propose a novel lip-to-speech generative adversarial network, Visual Context Attentional GAN (VCA-GAN), which can jointly model local and global lip movements during speech synthesis. Specifically, the proposed VCA-GAN synthesizes the speech from local lip visual features by finding a mapping function of viseme-to-phoneme, while global visual context is embedded into the intermediate layers of the generator to clarify the ambiguity in the mapping induced by homophene. To achieve this, a visual context attention module is proposed where it encodes global representations from the local visual features, and provides the desired global visual context corresponding to the given coarse speech

representation to the generator through audio-visual attention. In addition to the explicit modelling of local and global visual representations, synchronization learning is introduced as a form of contrastive learning that guides the generator to synthesize a speech in sync with the given input lip movements. Extensive experiments demonstrate that the proposed VCA-GAN outperforms existing state-of-the-art and is able to effectively synthesize the speech from multi-speaker that has been barely handled in the previous works.

[Non-convex Distributionally Robust Optimization: Non-asymptotic Analysis](#)

- Jikai Jin, Bohang Zhang, Haiyang Wang, Liwei Wang
- abstract@[open-review](#): Distributionally robust optimization (DRO) is a widely-used approach to learn models that are robust against distribution shift. Compared with the standard optimization setting, the objective function in DRO is more difficult to optimize, and most of the existing theoretical results make strong assumptions on the loss function. In this work we bridge the gap by studying DRO algorithms for general smooth non-convex losses. By carefully exploiting the specific form of the DRO objective, we are able to provide non-asymptotic convergence guarantees even though the objective function is possibly non-convex, non-smooth and has unbounded gradient noise. In particular, we prove that a special algorithm called the mini-batch normalized gradient descent with momentum, can find an $\mathcal{O}(\epsilon^{-4})$ gradient complexity. We also discuss the conditional value-at-risk (CVaR) setting, where we propose a penalized DRO objective based on a smoothed version of the CVaR that allows us to obtain a similar convergence guarantee. We finally verify our theoretical results in a number of tasks and find that the proposed algorithm can consistently achieve prominent acceleration.

[Goal-Aware Cross-Entropy for Multi-Target Reinforcement Learning](#)

- Kibom Kim, Min Whoo Lee, Yoonsung Kim, JeHwan Ryu, Minsu Lee, Byoung-Tak Zhang
- abstract@[open-review](#): Learning in a multi-target environment without prior knowledge about the targets requires a large amount of samples and makes generalization difficult. To solve this problem, it is important to be able to discriminate targets through semantic understanding. In this paper, we propose goal-aware cross-entropy (GACE) loss, that can be utilized in a self-supervised way using auto-labeled goal states alongside reinforcement learning. Based on the loss, we then devise goal-discriminative attention networks (GDAN) which utilize the goal-relevant information to focus on the given instruction. We evaluate the proposed methods on visual navigation and robot arm manipulation tasks with multi-target environments and show that GDAN outperforms the state-of-the-art methods in terms of task success ratio, sample efficiency, and generalization. Additionally, qualitative analyses demonstrate that our proposed method can help the agent become aware of and focus on the given instruction clearly, promoting goal-directed behavior.

[Smooth Normalizing Flows](#)

- Jonas Köhler, Andreas Krämer, Frank Noe
- abstract@[open-review](#): Normalizing flows are a promising tool for modeling probability distributions in physical systems. While state-of-the-art flows accurately approximate distributions and energies, applications in physics additionally require smooth energies to compute forces and higher-order derivatives. Furthermore, such densities are often defined on non-trivial topologies. A recent example are Boltzmann Generators for generating 3D-structures of peptides and small proteins. These generative models leverage the space of internal coordinates (dihedrals, angles, and bonds), which is a product of hypertori and compact intervals. In this work, we introduce a class of smooth mixture transformations working on both compact intervals and hypertori. Mixture transformations employ root-finding methods to invert them in practice, which has so far prevented bi-directional flow training. To this end, we show that parameter gradients and forces of such inverses can be computed from forward evaluations via the inverse function theorem. We demonstrate two advantages of such smooth flows: they allow training by force matching to simulation data and can be used as potentials in molecular dynamics simulations.

[MetaAvatar: Learning Animatable Clothed Human Models from Few Depth Images](#)

- Shaofei Wang, Marko Mihajlovic, Qianli Ma, Andreas Geiger, Siyu Tang
- abstract@[open-review](#): In this paper, we aim to create generalizable and controllable neural signed distance fields (SDFs) that represent clothed humans from monocular depth observations. Recent advances in deep learning, especially neural implicit representations, have enabled human shape reconstruction and controllable avatar generation from different sensor inputs. However, to generate realistic cloth deformations from novel input poses, watertight meshes or dense full-body scans are usually needed as inputs. Furthermore, due to the difficulty of effectively modeling pose-dependent cloth deformations for diverse body shapes and cloth types, existing approaches resort to per-subject/cloth-type optimization from scratch, which is computationally expensive. In contrast, we propose an approach that can quickly generate realistic clothed human avatars, represented as controllable neural SDFs, given only monocular depth images. We achieve this by using meta-learning to learn an initialization of a hypernetwork that predicts the parameters of neural SDFs. The hypernetwork is conditioned on human poses and represents a clothed neural avatar that deforms non-rigidly according to the input poses. Meanwhile, it is meta-learned to effectively incorporate priors of diverse body shapes and cloth types and thus can be much faster to fine-tune, compared to models trained from scratch. We qualitatively and quantitatively show that our approach outperforms state-of-the-art approaches that require complete meshes as inputs while our approach requires only depth frames as inputs and runs orders of magnitudes faster. Furthermore, we demonstrate that our meta-learned hypernetwork is very robust, being the first to generate avatars with realistic dynamic cloth deformations given as few as 8 monocular depth frames.

[Distributed Principal Component Analysis with Limited Communication](#)

- Foivos Alimisis, Peter Davies, Bart Vandereycken, Dan Alistarh
- abstract@[open-review](#): We study efficient distributed algorithms for the fundamental problem of principal component analysis and leading eigenvector computation on the sphere, when the data are randomly distributed among a set of computational nodes. We propose a new quantized variant of Riemannian gradient descent to solve this problem, and prove that the algorithm converges with high probability under a set of necessary spherical-convexity properties. We give bounds on the number of bits transmitted by the algorithm under common initialization schemes, and investigate the dependency on the problem dimension in each case.

[Newton-LESS: Sparsification without Trade-offs for the Sketched Newton Update](#)

- Michal Derezhinski, Jonathan Lacotte, Mert Pilanci, Michael W. Mahoney
- abstract@[open-review](#): In second-order optimization, a potential bottleneck can be computing the Hessian matrix of the optimized function at every iteration. Randomized sketching has emerged as a powerful technique for constructing estimates of the Hessian which can be used to perform approximate Newton steps. This involves multiplication by a random sketching matrix, which introduces a trade-off between the computational cost of sketching and the convergence rate of the optimization. A theoretically desirable but practically much too expensive choice is to use a dense Gaussian sketching matrix, which produces unbiased estimates of the exact Newton step and offers strong problem-independent convergence guarantees. We show that the Gaussian matrix can be drastically sparsified, substantially reducing the computational cost, without affecting its convergence properties in any way. This approach, called Newton-LESS, is based on a recently introduced sketching technique: LEverage Score Sparsified (LESS) embeddings. We prove that Newton-LESS enjoys nearly the same problem-independent local convergence rate as Gaussian embeddings for a large class of functions. In particular, this leads to a new state-of-the-art convergence result for an iterative least squares solver. Finally, we substantially extend LESS embeddings to include uniformly sparsified random sign matrices which can be implemented efficiently and perform well in numerical experiments.

Confident Anchor-Induced Multi-Source Free Domain Adaptation

- Jiahua Dong, Zhen Fang, Anjin Liu, Gan Sun, Tongliang Liu
- abstract@[open-review](#): Unsupervised domain adaptation has attracted appealing academic attentions by transferring knowledge from labeled source domain to unlabeled target domain. However, most existing methods assume the source data are drawn from a single domain, which cannot be successfully applied to explore complementarily transferable knowledge from multiple source domains with large distribution discrepancies. Moreover, they require access to source data during training, which are inefficient and unpractical due to privacy preservation and memory storage. To address these challenges, we develop a novel Confident-Anchor-induced multi-source-free Domain Adaptation (CAiDA) model, which is a pioneer exploration of knowledge adaptation from multiple source domains to the unlabeled target domain without any source data, but with only pre-trained source models. Specifically, a source-specific transferable perception module is proposed to automatically quantify the contributions of the complementary knowledge transferred from multi-source domains to the target domain. To generate pseudo labels for the target domain without access to the source data, we develop a confident-anchor-induced pseudo label generator by constructing a confident anchor group and assigning each unconfident target sample with a semantic-nearest confident anchor. Furthermore, a class-relationship-aware consistency loss is proposed to preserve consistent inter-class relationships by aligning soft confusion matrices across domains. Theoretical analysis answers why multi-source domains are better than a single source domain, and establishes a novel learning bound to show the effectiveness of exploiting multi-source domains. Experiments on several representative datasets illustrate the superiority of our proposed CAiDA model. The code is available at <https://github.com/Learning-group123/CAiDA>.

Word2Fun: Modelling Words as Functions for Diachronic Word Representation

- Benyou Wang, Emanuele Di Buccio, Massimo Melucci
- abstract@[open-review](#): Word meaning may change over time as a reflection of changes in human society. Therefore, modeling time in word representation is necessary for some diachronic tasks. Most existing diachronic word representation approaches train the embeddings separately for each pre-grouped time-stamped corpus and align these embeddings, e.g., by orthogonal projections, vector initialization, temporal referencing, and compass. However, not only does word meaning change in a short time, word meaning may also be subject to evolution over long timespans, thus resulting in a unified continuous process. A recent approach called DiffTime' models semantic evolution as functions parameterized by multiple-layer nonlinear neural networks over time. In this paper, we will carry on this line of work by learning explicit functions over time for each word. Our approach, calledWord2Fun', reduces the space complexity from $\mathcal{O}(TVD)$ to $\mathcal{O}(kVD)$ where k is a small constant ($k \ll T$). In particular, a specific instance based on polynomial functions could provably approximate any function modeling word evolution with a given negligible error thanks to the Weierstrass Approximation Theorem. The effectiveness of the proposed approach is evaluated in diverse tasks including time-aware word clustering, temporal analogy, and semantic change detection. Code at: <https://github.com/wabyking/Word2Fun.git>.

Iteratively Reweighted Least Squares for Basis Pursuit with Global Linear Convergence Rate

- Christian Kümmerle, Claudio Mayrink Verdun, Dominik Stöger
- abstract@[open-review](#): The recovery of sparse data is at the core of many applications in machine learning and signal processing. While such problems can be tackled using ℓ_1 -regularization as in the LASSO estimator and in the Basis Pursuit approach, specialized algorithms are typically required to solve the corresponding high-dimensional non-smooth optimization for large instances. Iteratively Reweighted Least Squares (IRLS) is a widely used algorithm for this purpose due to its excellent numerical performance. However, while existing theory is able to guarantee convergence of this algorithm to the minimizer, it does not provide a global convergence rate. In this paper, we prove that a variant of IRLS converges with a global linear rate to a sparse solution, i.e., with a linear error decrease occurring immediately from any initialization if the measurements fulfill the usual null space property assumption. We support our theory by numerical experiments showing that our linear rate captures the correct dimension dependence. We anticipate that our theoretical findings will lead to new insights for many other use cases of the IRLS algorithm, such as in low-rank matrix recovery.

Low-Rank Constraints for Fast Inference in Structured Models

- Justin Chiu, Yuntian Deng, Alexander Rush
- abstract@[open-review](#): Structured distributions, i.e. distributions over combinatorial spaces, are commonly used to learn latent probabilistic representations from observed data. However, scaling these models is bottlenecked by the high computational and memory complexity with respect to the size of the latent representations. Common models such as Hidden Markov Models (HMMs) and Probabilistic Context-Free Grammars (PCFGs) require time and space quadratic and cubic in the number of hidden states respectively. This work demonstrates a simple approach to reduce the computational and memory complexity of a large class of structured models. We show that by viewing the central inference step as a matrix-vector product and using a low-rank constraint, we can trade off model expressivity and speed via the rank. Experiments with neural parameterized structured models for language modeling, polyphonic music modeling, unsupervised grammar induction, and video modeling show that our approach matches the accuracy of standard models at large state spaces while providing practical speedups.

Accumulative Poisoning Attacks on Real-time Data

- Tianyu Pang, Xiao Yang, Yinpeng Dong, Hang Su, Jun Zhu
- abstract@[open-review](#): Collecting training data from untrusted sources exposes machine learning services to poisoning adversaries, who maliciously manipulate training data to degrade the model accuracy. When trained on offline datasets, poisoning adversaries have to inject the poisoned data in advance before training, and the order of feeding these poisoned batches into the model is stochastic. In contrast, practical systems are more usually trained/fine-tuned on sequentially captured real-time data, in which case poisoning adversaries could dynamically poison each data batch according to the current model state. In this paper, we focus on the real-time settings and propose a new attacking strategy, which affiliates an accumulative phase with poisoning attacks to secretly (i.e., without affecting accuracy) magnify the destructive effect of a (poisoned) trigger batch. By mimicking online learning and federated learning on MNIST and CIFAR-10, we show that model accuracy significantly drops by a single update step on the trigger batch after the accumulative phase. Our work validates that a well-designed but straightforward attacking strategy can dramatically amplify the poisoning effects, with no need to explore complex techniques.

UCB-based Algorithms for Multinomial Logistic Regression Bandits

- Sanae Amani, Christos Thrampoulidis
- abstract@[open-review](#): Out of the rich family of generalized linear bandits, perhaps the most well studied ones are logistic bandits that are used in problems with binary rewards: for instance, when the learner aims to maximize the profit over a user that can select one of two possible outcomes (e.g., click' vsno-click'). Despite remarkable recent progress and improved algorithms for logistic bandits, existing works do not address practical situations where the number of outcomes that can be selected by the user is larger than two (e.g., click', show me later', never show again', no click'). In this paper, we study such an extension. We use multinomial logit (MNL) to model the probability of each one of $K+1$ possible outcomes (+1 stands for the 'not click' outcome): we assume that for a learner's action x_t , the user selects one of $K+1$ outcomes, say outcome i , with a MNL probabilistic model with corresponding unknown parameter $\bar{\theta}_i$. Each outcome i is also associated with a revenue parameter ρ_i and the goal is to maximize the expected revenue. For this problem, we present MNL-UCB, an upper confidence bound (UCB)-based algorithm, that achieves regret $\tilde{\mathcal{O}}(dK\sqrt{T})$ with small dependency on problem-dependent constants that can otherwise be arbitrarily large and lead to loose regret bounds. We present numerical simulations that corroborate our theoretical results.

Estimating the Long-Term Effects of Novel Treatments

- Keith Battocchi, Eleanor Dillon, Maggie Hei, Greg Lewis, Miruna Oprescu, Vasilis Syrgkanis
- abstract@[open-review](#): Policy makers often need to estimate the long-term effects of novel treatments, while only having historical data of older treatment options. We propose a surrogate-based approach using a long-term dataset where only past treatments were administered and a short-term dataset where novel treatments have been administered. Our approach generalizes previous surrogate-style methods, allowing for continuous treatments and serially-correlated treatment policies while maintaining consistency and root-n asymptotically normal estimates under a Markovian assumption on the data and the observational policy. Using a semi-synthetic dataset on customer incentives from a major corporation, we evaluate the performance of our method and discuss solutions to practical challenges when deploying our methodology.

Dual Progressive Prototype Network for Generalized Zero-Shot Learning

- Chaoqun Wang, Shaobo Min, Xuejin Chen, Xiaoyan Sun, Houqiang Li
- abstract@[open-review](#): Generalized Zero-Shot Learning (GZSL) aims to recognize new categories with auxiliary semantic information, e.g., category attributes. In this paper, we handle the critical issue of domain shift problem, i.e., confusion between seen and unseen categories, by progressively improving cross-domain transferability and category discriminability of visual representations. Our approach, named Dual Progressive Prototype Network (DPPN), constructs two types of prototypes that record prototypical visual patterns for attributes and categories, respectively. With attribute prototypes, DPPN alternately searches attribute-related local regions and updates corresponding attribute prototypes to progressively explore accurate attribute-region correspondence. This enables DPPN to produce visual representations with accurate attribute localization ability, which benefits the semantic-visual alignment and representation transferability. Besides, along with progressive attribute localization, DPPN further projects category prototypes into multiple spaces to progressively repel visual representations from different categories, which boosts category discriminability. Both attribute and category prototypes are collaboratively learned in a unified framework, which makes visual representations of DPPN transferable and distinctive. Experiments on four benchmarks prove that DPPN effectively alleviates the domain shift problem in GZSL.

Derivative-Free Policy Optimization for Linear Risk-Sensitive and Robust Control Design: Implicit Regularization and Sample Complexity

- Kaiqing Zhang, Xiangyuan Zhang, Bin Hu, Tamer Basar
- abstract@[open-review](#): Direct policy search serves as one of the workhorses in modern reinforcement learning (RL), and its applications in continuous control tasks have recently attracted increasing attention. In this work, we investigate the convergence theory of policy gradient (PG) methods for learning the linear risk-sensitive and robust controller. In particular, we develop PG methods that can be implemented in a derivative-free fashion by sampling system trajectories, and establish both global convergence and sample complexity results in the solutions of two fundamental settings in risk-sensitive and robust control: the finite-horizon linear exponential quadratic Gaussian, and the finite-horizon linear-quadratic disturbance attenuation problems. As a by-product, our results also provide the first sample complexity for the global convergence of PG methods on solving zero-sum linear-quadratic dynamic games, a nonconvex-nonconcave minimax optimization problem that serves as a baseline setting in multi-agent reinforcement learning (MARL) with continuous spaces. One feature of our algorithms is that during the learning phase, a certain level of robustness/risk-sensitivity of the controller is preserved, which we termed as the implicit regularization property, and is an essential requirement in safety-critical control systems.

G-PATE: Scalable Differentially Private Data Generator via Private Aggregation of Teacher Discriminators

- Yunhui Long, Boxin Wang, Zhuolin Yang, Bhavya Kailkhura, Aston Zhang, Carl Gunter, Bo Li
- abstract@[open-review](#): Recent advances in machine learning have largely benefited from the massive accessible training data. However, large-scale data sharing has raised great privacy concerns. In this work, we propose a novel privacy-preserving data Generative model based on the PATE framework (G-PATE), aiming to train a scalable differentially private data generator that preserves high generated data utility. Our approach leverages generative adversarial nets to generate data, combined with private aggregation among different discriminators to ensure strong privacy guarantees. Compared to existing approaches, G-PATE significantly improves the use of privacy budgets. In particular, we train a student data generator with an ensemble of teacher discriminators and propose a novel private gradient aggregation mechanism to ensure differential privacy on all information that flows from teacher discriminators to the student generator. In addition, with random projection and gradient discretization, the proposed gradient aggregation mechanism is able to effectively deal with high-dimensional gradient vectors. Theoretically, we prove that G-PATE ensures differential privacy for the data generator. Empirically, we demonstrate the superiority of G-PATE over prior work through extensive experiments. We show that G-PATE is the first work being able to generate high-dimensional image data with high data utility under limited privacy budgets ($\$\\backslash varepsilon \\leq 1\$$). Our code is available at <https://github.com/AI-secure/G-PATE>.

On the Existence of The Adversarial Bayes Classifier

- Pranjal Awasthi, Natalie Frank, Mehryar Mohri
- abstract@[open-review](#): Adversarial robustness is a critical property in a variety of modern machine learning applications. While it has been the subject of several recent theoretical studies, many important questions related to adversarial robustness are still open. In this work, we study a fundamental question regarding Bayes optimality for adversarial robustness. We provide general sufficient conditions under which the existence of a Bayes optimal classifier can be guaranteed for adversarial robustness. Our results can provide a useful tool for a subsequent study of surrogate losses in adversarial robustness and their consistency properties.

Convex-Concave Min-Max Stackelberg Games

- Denizalp Goktas, Amy Greenwald
- abstract@[open-review](#): Min-max optimization problems (i.e., min-max games) have been attracting a great deal of attention because of their applicability to a wide range of machine learning problems. Although significant progress has been made recently, the literature to date has focused on games with independent strategy sets; little is known about solving games with dependent strategy sets, which can be characterized as min-max Stackelberg games. We introduce two first-order methods that solve a large class of convex-concave min-max Stackelberg games, and show that our methods converge in polynomial time. Min-max Stackelberg games were first studied by Wald, under the posthumous name of Wald's maximin model, a variant of which is the main paradigm used in robust optimization, which means that our methods can likewise solve many convex robust optimization problems. We observe that the computation of competitive equilibria in Fisher markets also comprises a min-max Stackelberg game. Further, we demonstrate the efficacy and efficiency of our algorithms in practice by computing competitive equilibria in Fisher markets with varying utility structures. Our experiments suggest potential ways to extend our theoretical results, by demonstrating how different smoothness properties can affect the convergence rate of our algorithms.

Misspecified Gaussian Process Bandit Optimization

- Ilija Bogunovic, Andreas Krause
- abstract@[open-review](#): We consider the problem of optimizing a black-box function based on noisy bandit feedback. Kernelized bandit algorithms have shown strong empirical and theoretical performance for this problem. They heavily rely on the assumption that the model is well-specified, however, and can fail without it. Instead, we introduce and address a $\backslash emph{misspecified}$ kernelized bandit setting where the unknown function can be $\backslash epsilon--$

uniformly approximated by a function with a bounded norm in some Reproducing Kernel Hilbert Space (RKHS). We design efficient and practical algorithms whose performance degrades minimally in the presence of model misspecification. Specifically, we present two algorithms based on Gaussian process (GP) methods: an optimistic EC-GP-UCB algorithm that requires knowing the misspecification error, and Phased GP Uncertainty Sampling, an elimination-type algorithm that can adapt to unknown model misspecification. We provide upper bounds on their cumulative regret in terms of $\$\\epsilon$, the time horizon, and the underlying kernel, and we show that our algorithm achieves optimal dependence on $\$\\epsilon$ with no prior knowledge of misspecification. In addition, in a stochastic contextual setting, we show that EC-GP-UCB can be effectively combined with the regret bound balancing strategy and attain similar regret bounds despite not knowing $\$\\epsilon$.

[Visual Adversarial Imitation Learning using Variational Models](#)

- Rafael Rafailov, Tianhe Yu, Aravind Rajeswaran, Chelsea Finn
- abstract@[open-review](#): Reward function specification, which requires considerable human effort and iteration, remains a major impediment for learning behaviors through deep reinforcement learning. In contrast, providing visual demonstrations of desired behaviors presents an easier and more natural way to teach agents. We consider a setting where an agent is provided a fixed dataset of visual demonstrations illustrating how to perform a task, and must learn to solve the task using the provided demonstrations and unsupervised environment interactions. This setting presents a number of challenges including representation learning for visual observations, sample complexity due to high dimensional spaces, and learning instability due to the lack of a fixed reward or learning signal. Towards addressing these challenges, we develop a variational model-based adversarial imitation learning (V-MAIL) algorithm. The model-based approach provides a strong signal for representation learning, enables sample efficiency, and improves the stability of adversarial training by enabling on-policy learning. Through experiments involving several vision-based locomotion and manipulation tasks, we find that V-MAIL learns successful visuomotor policies in a sample-efficient manner, has better stability compared to prior work, and also achieves higher asymptotic performance. We further find that by transferring the learned models, V-MAIL can learn new tasks from visual demonstrations without any additional environment interactions. All results including videos can be found online at <https://sites.google.com/view/variational-mail>

[Object-Aware Regularization for Addressing Causal Confusion in Imitation Learning](#)

- Jongjin Park, Younggyo Seo, Chang Liu, Li Zhao, Tao Qin, Jinwoo Shin, Tie-Yan Liu
- abstract@[open-review](#): Behavioral cloning has proven to be effective for learning sequential decision-making policies from expert demonstrations. However, behavioral cloning often suffers from the causal confusion problem where a policy relies on the noticeable effect of expert actions due to the strong correlation but not the cause we desire. This paper presents Object-aware REgularizatiOn (OREO), a simple technique that regularizes an imitation policy in an object-aware manner. Our main idea is to encourage a policy to uniformly attend to all semantic objects, in order to prevent the policy from exploiting nuisance variables strongly correlated with expert actions. To this end, we introduce a two-stage approach: (a) we extract semantic objects from images by utilizing discrete codes from a vector-quantized variational autoencoder, and (b) we randomly drop the units that share the same discrete code together, i.e., masking out semantic objects. Our experiments demonstrate that OREO significantly improves the performance of behavioral cloning, outperforming various other regularization and causality-based methods on a variety of Atari environments and a self-driving CARLA environment. We also show that our method even outperforms inverse reinforcement learning methods trained with a considerable amount of environment interaction.

[Reliable and Trustworthy Machine Learning for Health Using Dataset Shift Detection](#)

- Chunjong Park, Anas Awadalla, Tadayoshi Kohno, Shwetak Patel
- abstract@[open-review](#): Unpredictable ML model behavior on unseen data, especially in the health domain, raises serious concerns about its safety as repercussions for mistakes can be fatal. In this paper, we explore the feasibility of using state-of-the-art out-of-distribution detectors for reliable and trustworthy diagnostic predictions. We select publicly available deep learning models relating to various health conditions (e.g., skin cancer, lung sound, and Parkinson's disease) using various input data types (e.g., image, audio, and motion data). We demonstrate that these models show unreasonable predictions on out-of-distribution datasets. We show that Mahalanobis distance- and Gram matrices-based out-of-distribution detection methods are able to detect out-of-distribution data with high accuracy for the health models that operate on different modalities. We then translate the out-of-distribution score into a human interpretable \textsc{confidence score} to investigate its effect on the users' interaction with health ML applications. Our user study shows that the \textsc{confidence score} helped the participants only trust the results with a high score to make a medical decision and disregard results with a low score. Through this work, we demonstrate that dataset shift is a critical piece of information for high-stake ML applications, such as medical diagnosis and healthcare, to provide reliable and trustworthy predictions to the users.

[Multiclass Boosting and the Cost of Weak Learning](#)

- Nataly Brukhim, Elad Hazan, Shay Moran, Indraneel Mukherjee, Robert E. Schapire
- abstract@[open-review](#): Boosting is an algorithmic approach which is based on the idea of combining weak and moderately inaccurate hypotheses to a strong and accurate one. In this work we study multiclass boosting with a possibly large number of classes or categories. Multiclass boosting can be formulated in various ways. Here, we focus on an especially natural formulation in which the weak hypotheses are assumed to belong to an "easy-to-learn" base class, and the weak learner is an agnostic PAC learner for that class with respect to the standard classification loss. This is in contrast with other, more complicated losses as have often been considered in the past. The goal of the overall boosting algorithm is then to learn a combination of weak hypotheses by repeatedly calling the weak learner. We study the resources required for boosting, especially how they depend on the number of classes k , for both the booster and weak learner. We find that the boosting algorithm itself only requires $O(\log k)$ samples, as we show by analyzing a variant of AdaBoost for our setting. In stark contrast, assuming typical limits on the number of weak-learner calls, we prove that the number of samples required by a weak learner is at least polynomial in k , exponentially more than the number of samples needed by the booster. Alternatively, we prove that the weak learner's accuracy parameter must be smaller than an inverse polynomial in k , showing that the returned weak hypotheses must be nearly the best in their class when k is large. We also prove a trade-off between number of oracle calls and the resources required of the weak learner, meaning that the fewer calls to the weak learner the more that is demanded on each call.

[Partition-Based Formulations for Mixed-Integer Optimization of Trained ReLU Neural Networks](#)

- Calvin Tsay, Jan Kronqvist, Alexander Thebelt, Ruth Misener
- abstract@[open-review](#): This paper introduces a class of mixed-integer formulations for trained ReLU neural networks. The approach balances model size and tightness by partitioning node inputs into a number of groups and forming the convex hull over the partitions via disjunctive programming. At one extreme, one partition per input recovers the convex hull of a node, i.e., the tightest possible formulation for each node. For fewer partitions, we develop smaller relaxations that approximate the convex hull, and show that they outperform existing formulations. Specifically, we propose strategies for partitioning variables based on theoretical motivations and validate these strategies using extensive computational experiments. Furthermore, the proposed scheme complements known algorithmic approaches, e.g., optimization-based bound tightening captures dependencies within a partition.

[Hyperparameter Optimization Is Deceiving Us, and How to Stop It](#)

- A. Feder Cooper, Yucheng Lu, Jessica Forde, Christopher M. De Sa
- abstract@[open-review](#): Recent empirical work shows that inconsistent results based on choice of hyperparameter optimization (HPO) configuration are a widespread problem in ML research. When comparing two algorithms J and K searching one subspace can yield the conclusion that J outperforms K, whereas searching another can entail the opposite. In short, the way we choose hyperparameters can deceive us. We provide a theoretical complement to

this prior work, arguing that, to avoid such deception, the process of drawing conclusions from HPO should be made more rigorous. We call this process epistemic hyperparameter optimization (EHPO), and put forth a logical framework to capture its semantics and how it can lead to inconsistent conclusions about performance. Our framework enables us to prove EHPO methods that are guaranteed to be defended against deception, given bounded compute time budget t . We demonstrate our framework's utility by proving and empirically validating a defended variant of random search.

[On the Convergence Theory of Debiased Model-Agnostic Meta-Reinforcement Learning](#)

- Alireza Fallah, Kristian Georgiev, Aryan Mokhtari, Asuman Ozdaglar
- abstract@[open-review](#): We consider Model-Agnostic Meta-Learning (MAML) methods for Reinforcement Learning (RL) problems, where the goal is to find a policy using data from several tasks represented by Markov Decision Processes (MDPs) that can be updated by one step of stochastic policy gradient for the realized MDP. In particular, using stochastic gradients in MAML update steps is crucial for RL problems since computation of exact gradients requires access to a large number of possible trajectories. For this formulation, we propose a variant of the MAML method, named Stochastic Gradient Meta-Reinforcement Learning (SG-MRL), and study its convergence properties. We derive the iteration and sample complexity of SG-MRL to find an ϵ -first-order stationary point, which, to the best of our knowledge, provides the first convergence guarantee for model-agnostic meta-reinforcement learning algorithms. We further show how our results extend to the case where more than one step of stochastic policy gradient method is used at test time. Finally, we empirically compare SG-MRL and MAML in several deep RL environments.

[3D Pose Transfer with Correspondence Learning and Mesh Refinement](#)

- Chaoyue Song, Jiacheng Wei, Ruibo Li, Fayao Liu, Guosheng Lin
- abstract@[open-review](#): 3D pose transfer is one of the most challenging 3D generation tasks. It aims to transfer the pose of a source mesh to a target mesh and keep the identity (e.g., body shape) of the target mesh. Some previous works require key point annotations to build reliable correspondence between the source and target meshes, while other methods do not consider any shape correspondence between sources and targets, which leads to limited generation quality. In this work, we propose a correspondence-refinement network to achieve the 3D pose transfer for both human and animal meshes. The correspondence between source and target meshes is first established by solving an optimal transport problem. Then, we warp the source mesh according to the dense correspondence and obtain a coarse warped mesh. The warped mesh will be better refined with our proposed Elastic Instance Normalization, which is a conditional normalization layer and can help to generate high-quality meshes. Extensive experimental results show that the proposed architecture can effectively transfer the poses from source to target meshes and produce better results with satisfied visual performance than state-of-the-art methods.

[Framing RNN as a kernel method: A neural ODE approach](#)

- Adeline Fermanian, Pierre Marion, Jean-Philippe Vert, Gérard Biau
- abstract@[open-review](#): Building on the interpretation of a recurrent neural network (RNN) as a continuous-time neural differential equation, we show, under appropriate conditions, that the solution of a RNN can be viewed as a linear function of a specific feature set of the input sequence, known as the signature. This connection allows us to frame a RNN as a kernel method in a suitable reproducing kernel Hilbert space. As a consequence, we obtain theoretical guarantees on generalization and stability for a large class of recurrent networks. Our results are illustrated on simulated datasets.

[Contextual Similarity Aggregation with Self-attention for Visual Re-ranking](#)

- Jianbo Ouyang, Hui Wu, Min Wang, Wengang Zhou, Houqiang Li
- abstract@[open-review](#): In content-based image retrieval, the first-round retrieval result by simple visual feature comparison may be unsatisfactory, which can be refined by visual re-ranking techniques. In image retrieval, it is observed that the contextual similarity among the top-ranked images is an important clue to distinguish the semantic relevance. Inspired by this observation, in this paper, we propose a visual re-ranking method by contextual similarity aggregation with self-attention. In our approach, for each image in the top-K ranking list, we represent it into an affinity feature vector by comparing it with a set of anchor images. Then, the affinity features of the top-K images are refined by aggregating the contextual information with a transformer encoder. Finally, the affinity features are used to recalculate the similarity scores between the query and the top-K images for re-ranking of the latter. To further improve the robustness of our re-ranking model and enhance the performance of our method, a new data augmentation scheme is designed. Since our re-ranking model is not directly involved with the visual feature used in the initial retrieval, it is ready to be applied to retrieval result lists obtained from various retrieval algorithms. We conduct comprehensive experiments on four benchmark datasets to demonstrate the generality and effectiveness of our proposed visual re-ranking method.

[Can Information Flows Suggest Targets for Interventions in Neural Circuits?](#)

- Praveen Venkatesh, Sanghamitra Dutta, Neil Mehta, Pulkit Grover
- abstract@[open-review](#): Motivated by neuroscientific and clinical applications, we empirically examine whether observational measures of information flow can suggest interventions. We do so by performing experiments on artificial neural networks in the context of fairness in machine learning, where the goal is to induce fairness in the system through interventions. Using our recently developed M-information flow framework, we measure the flow of information about the true label (responsible for accuracy, and hence desirable), and separately, the flow of information about a protected attribute (responsible for bias, and hence undesirable) on the edges of a trained neural network. We then compare the flow magnitudes against the effect of intervening on those edges by pruning. We show that pruning edges that carry larger information flows about the protected attribute reduces bias at the output to a greater extent. This demonstrates that M-information flow can meaningfully suggest targets for interventions, answering the title's question in the affirmative. We also evaluate bias-accuracy tradeoffs for different intervention strategies, to analyze how one might use estimates of desirable and undesirable information flows (here, accuracy and bias flows) to inform interventions that preserve the former while reducing the latter.

[AutoBalance: Optimized Loss Functions for Imbalanced Data](#)

- Mingchen Li, Xuechen Zhang, Christos Thrampoulidis, Jiasi Chen, Samet Oymak
- abstract@[open-review](#): Imbalanced datasets are commonplace in modern machine learning problems. The presence of under-represented classes or groups with sensitive attributes results in concerns about generalization and fairness. Such concerns are further exacerbated by the fact that large capacity deep nets can perfectly fit the training data and appear to achieve perfect accuracy and fairness during training, but perform poorly during test. To address these challenges, we propose AutoBalance, a bi-level optimization framework that automatically designs a training loss function to optimize a blend of accuracy and fairness-seeking objectives. Specifically, a lower-level problem trains the model weights, and an upper-level problem tunes the loss function by monitoring and optimizing the desired objective over the validation data. Our loss design enables personalized treatment for classes/groups by employing a parametric cross-entropy loss and individualized data augmentation schemes. We evaluate the benefits and performance of our approach for the application scenarios of imbalanced and group-sensitive classification. Extensive empirical evaluations demonstrate the benefits of AutoBalance over state-of-the-art approaches. Our experimental findings are complemented with theoretical insights on loss function design and the benefits of the train-validation split. All code is available open-source.

[SyncTwin: Treatment Effect Estimation with Longitudinal Outcomes](#)

- Zhaozhi Qian, Yao Zhang, Ioana Bica, Angela Wood, Mihaela van der Schaar
- abstract@[open-review](#): Most of the medical observational studies estimate the causal treatment effects using electronic health records (EHR), where a patient's covariates and outcomes are both observed longitudinally. However, previous methods focus only on adjusting for the covariates while neglecting the temporal structure in the outcomes. To bridge the gap, this paper develops a new method, SyncTwin, that learns a patient-specific time-constant representation from the pre-treatment observations. SyncTwin issues counterfactual prediction of a target patient by constructing a synthetic twin that closely matches the target in representation. The reliability of the estimated treatment effect can be assessed by comparing the observed and synthetic pre-treatment outcomes. The medical experts can interpret the estimate by examining the most important contributing individuals to the synthetic twin. In the real-data experiment, SyncTwin successfully reproduced the findings of a randomized controlled clinical trial using observational data, which demonstrates its usability in the complex real-world EHR.

[Statistical Query Lower Bounds for List-Decodable Linear Regression](#)

- Ilias Diakonikolas, Daniel Kane, Ankit Pensia, Thanasis Pittas, Alistair Stewart
- abstract@[open-review](#): We study the problem of list-decodable linear regression, where an adversary can corrupt a majority of the examples. Specifically, we are given a set T of labeled examples $(x, y) \in \mathbb{R}^d \times \mathbb{R}$ and a parameter $0 < \alpha < 1/2$ such that an α -fraction of the points in T are i.i.d. samples from a linear regression model with Gaussian covariates, and the remaining $(1-\alpha)$ -fraction of the points are drawn from an arbitrary noise distribution. The goal is to output a small list of hypothesis vectors such that at least one of them is close to the target regression vector. Our main result is a Statistical Query (SQ) lower bound of $d^{\mathrm{poly}(1/\alpha)}$ for this problem. Our SQ lower bound qualitatively matches the performance of previously developed algorithms, providing evidence that current upper bounds for this task are nearly best possible.

[Unsupervised Motion Representation Learning with Capsule Autoencoders](#)

- Ziwei Xu, Xudong Shen, Yongkang Wong, Mohan S. Kankanhalli
- abstract@[open-review](#): We propose the Motion Capsule Autoencoder (MCAE), which addresses a key challenge in the unsupervised learning of motion representations: transformation invariance. MCAE models motion in a two-level hierarchy. In the lower level, a spatio-temporal motion signal is divided into short, local, and semantic-agnostic snippets. In the higher level, the snippets are aggregated to form full-length semantic-aware segments. For both levels, we represent motion with a set of learned transformation invariant templates and the corresponding geometric transformations by using capsule autoencoders of a novel design. This leads to a robust and efficient encoding of viewpoint changes. MCAE is evaluated on a novel Trajectory20 motion dataset and various real-world skeleton-based human action datasets. Notably, it achieves better results than baselines on Trajectory20 with considerably fewer parameters and state-of-the-art performance on the unsupervised skeleton-based action recognition task.

[VigDet: Knowledge Informed Neural Temporal Point Process for Coordination Detection on Social Media](#)

- Yizhou Zhang, Karishma Sharma, Yan Liu
- abstract@[open-review](#): Recent years have witnessed an increasing use of coordinated accounts on social media, operated by misinformation campaigns to influence public opinion and manipulate social outcomes. Consequently, there is an urgent need to develop an effective methodology for coordinated group detection to combat the misinformation on social media. However, existing works suffer from various drawbacks, such as, either limited performance due to extreme reliance on predefined signatures of coordination, or instead an inability to address the natural sparsity of account activities on social media with useful prior domain knowledge. Therefore, in this paper, we propose a coordination detection framework incorporating neural temporal point process with prior knowledge such as temporal logic or pre-defined filtering functions. Specifically, when modeling the observed data from social media with neural temporal point process, we jointly learn a Gibbs-like distribution of group assignment based on how consistent an assignment is to (1) the account embedding space and (2) the prior knowledge. To address the challenge that the distribution is hard to be efficiently computed and sampled from, we design a theoretically guaranteed variational inference approach to learn a mean-field approximation for it. Experimental results on a real-world dataset show the effectiveness of our proposed method compared to the SOTA model in both unsupervised and semi-supervised settings. We further apply our model on a COVID-19 Vaccine Tweets dataset. The detection result suggests the presence of suspicious coordinated efforts on spreading misinformation about COVID-19 vaccines.

[An Improved Analysis and Rates for Variance Reduction under Without-replacement Sampling Orders](#)

- Xinmeng Huang, Kun Yuan, Xianghui Mao, Wotao Yin
- abstract@[open-review](#): When applying a stochastic algorithm, one must choose an order to draw samples. The practical choices are without-replacement sampling orders, which are empirically faster and more cache-friendly than uniform-iid-sampling but often have inferior theoretical guarantees. Without-replacement sampling is well understood only for SGD without variance reduction. In this paper, we will improve the convergence analysis and rates of variance reduction under without-replacement sampling orders for composite finite-sum minimization. Our results are in two-folds. First, we develop a damped variant of Finito called Prox-DFinito and establish its convergence rates with random reshuffling, cyclic sampling, and shuffling-once, under both generally and strongly convex scenarios. These rates match full-batch gradient descent and are state-of-the-art compared to the existing results for without-replacement sampling with variance-reduction. Second, our analysis can gauge how the cyclic order will influence the rate of cyclic sampling and, thus, allows us to derive the optimal fixed ordering. In the highly data-heterogeneous scenario, Prox-DFinito with optimal cyclic sampling can attain a sample-size-independent convergence rate, which, to our knowledge, is the first result that can match with uniform-iid-sampling with variance reduction. We also propose a practical method to discover the optimal cyclic ordering numerically.

[Exploring Forensic Dental Identification with Deep Learning](#)

- Yuan Liang, Weikun Han, Liang Qiu, Chen Wu, Yiting Shao, Kun Wang, Lei He
- abstract@[open-review](#): Dental forensic identification targets to identify persons with dental traces. The task is vital for the investigation of criminal scenes and mass disasters because of the resistance of dental structures and the wide-existence of dental imaging. However, no widely accepted automated solution is available for this labour-costly task. In this work, we pioneer to study deep learning for dental forensic identification based on panoramic radiographs. We construct a comprehensive benchmark with various dental variations that can adequately reflect the difficulties of the task. By considering the task's unique challenges, we propose FoID, a deep learning method featured by: (i) clinical-inspired attention localization, (ii) domain-specific augmentations that enable instance discriminative learning, and (iii) transformer-based self-attention mechanism that dynamically reasons the relative importance of attentions. We show that FoID can outperform traditional approaches by at least 22.98% in terms of Rank-1 accuracy, and outperform strong CNN baselines by at least 10.50% in terms of mean Average Precision (mAP). Moreover, extensive ablation studies verify the effectiveness of each building blocks of FoID. Our work can be a first step towards the automated system for forensic identification among large-scale multi-site databases. Also, the proposed techniques, e.g., self-attention mechanism, can also be meaningful for other identification tasks, e.g., pedestrian re-identification. Related data and codes can be found at <https://github.com/liangyuandg/FoID>.

[Learning to Generate Realistic Noisy Images via Pixel-level Noise-aware Adversarial Training](#)

- Yuanhao Cai, Xiaowan Hu, Haoqian Wang, Yulun Zhang, Hanspeter Pfister, Donglai Wei

- abstract@[open-review](#): Existing deep learning real denoising methods require a large amount of noisy-clean image pairs for supervision. Nonetheless, capturing a real noisy-clean dataset is an unacceptable expensive and cumbersome procedure. To alleviate this problem, this work investigates how to generate realistic noisy images. Firstly, we formulate a simple yet reasonable noise model that treats each real noisy pixel as a random variable. This model splits the noisy image generation problem into two sub-problems: image domain alignment and noise domain alignment. Subsequently, we propose a novel framework, namely Pixel-level Noise-aware Generative Adversarial Network (PNGAN). PNGAN employs a pre-trained real denoiser to map the fake and real noisy images into a nearly noise-free solution space to perform image domain alignment. Simultaneously, PNGAN establishes a pixel-level adversarial training to conduct noise domain alignment. Additionally, for better noise fitting, we present an efficient architecture Simple Multi-scale Network (SMNet) as the generator. Qualitative validation shows that noise generated by PNGAN is highly similar to real noise in terms of intensity and distribution. Quantitative experiments demonstrate that a series of denoisers trained with the generated noisy images achieve state-of-the-art (SOTA) results on four real denoising benchmarks.

[Multi-Agent Reinforcement Learning for Active Voltage Control on Power Distribution Networks](#)

- Jianhong Wang, Wangkun Xu, Yunjie Gu, Wenbin Song, Tim C Green
- abstract@[open-review](#): This paper presents a problem in power networks that creates an exciting and yet challenging real-world scenario for application of multi-agent reinforcement learning (MARL). The emerging trend of decarbonisation is placing excessive stress on power distribution networks. Active voltage control is seen as a promising solution to relieve power congestion and improve voltage quality without extra hardware investment, taking advantage of the controllable apparatuses in the network, such as roof-top photovoltaics (PVs) and static var compensators (SVCs). These controllable apparatuses appear in a vast number and are distributed in a wide geographic area, making MARL a natural candidate. This paper formulates the active voltage control problem in the framework of Dec-POMDP and establishes an open-source environment. It aims to bridge the gap between the power community and the MARL community and be a drive force towards real-world applications of MARL algorithms. Finally, we analyse the special characteristics of the active voltage control problems that cause challenges (e.g. interpretability) for state-of-the-art MARL approaches, and summarise the potential directions.

[Looking Beyond Single Images for Contrastive Semantic Segmentation Learning](#)

- FEIHU ZHANG, Philip Torr, Rene Ranftl, Stephan Richter
- abstract@[open-review](#): We present an approach to contrastive representation learning for semantic segmentation. Our approach leverages the representational power of existing feature extractors to find corresponding regions across images. These cross-image correspondences are used as auxiliary labels to guide the pixel-level selection of positive and negative samples for more effective contrastive learning in semantic segmentation. We show that auxiliary labels can be generated from a variety of feature extractors, ranging from image classification networks that have been trained using unsupervised contrastive learning to segmentation models that have been trained on a small amount of labeled data. We additionally introduce a novel metric for rapidly judging the quality of a given auxiliary-labeling strategy, and empirically analyze various factors that influence the performance of contrastive learning for semantic segmentation. We demonstrate the effectiveness of our method both in the low-data as well as the high-data regime on various datasets. Our experiments show that contrastive learning with our auxiliary-labeling approach consistently boosts semantic segmentation accuracy when compared to standard ImageNet pretraining and outperforms existing approaches of contrastive and semi-supervised semantic segmentation.

[A Constant Approximation Algorithm for Sequential Random-Order No-Substitution k-Median Clustering](#)

- Tom Hess, Michal Moshkovitz, Sivan Sabato
- abstract@[open-review](#): We study k-median clustering under the sequential no-substitution setting. In this setting, a data stream is sequentially observed, and some of the points are selected by the algorithm as cluster centers. However, a point can be selected as a center only immediately after it is observed, before observing the next point. In addition, a selected center cannot be substituted later. We give the first algorithm for this setting that obtains a constant approximation factor on the optimal cost under a random arrival order, an exponential improvement over previous work. This is also the first constant approximation guarantee that holds without any structural assumptions on the input data. Moreover, the number of selected centers is only quasi-linear in k. Our algorithm and analysis are based on a careful cost estimation that avoids outliers, a new concept of a linear bin division, and a multi-scale approach to center selection.

[Dangers of Bayesian Model Averaging under Covariate Shift](#)

- Pavel Izmailov, Patrick Nicholson, Sanae Lotfi, Andrew G. Wilson
- abstract@[open-review](#): Approximate Bayesian inference for neural networks is considered a robust alternative to standard training, often providing good performance on out-of-distribution data. However, Bayesian neural networks (BNNs) with high-fidelity approximate inference via full-batch Hamiltonian Monte Carlo achieve poor generalization under covariate shift, even underperforming classical estimation. We explain this surprising result, showing how a Bayesian model average can in fact be problematic under covariate shift, particularly in cases where linear dependencies in the input features cause a lack of posterior contraction. We additionally show why the same issue does not affect many approximate inference procedures, or classical maximum a-posteriori (MAP) training. Finally, we propose novel priors that improve the robustness of BNNs to many sources of covariate shift.

[Learning Equilibria in Matching Markets from Bandit Feedback](#)

- Meena Jagadeesan, Alexander Wei, Yixin Wang, Michael Jordan, Jacob Steinhardt
- abstract@[open-review](#): Large-scale, two-sided matching platforms must find market outcomes that align with user preferences while simultaneously learning these preferences from data. But since preferences are inherently uncertain during learning, the classical notion of stability (Gale and Shapley, 1962; Shapley and Shubik, 1971) is unattainable in these settings. To bridge this gap, we develop a framework and algorithms for learning stable market outcomes under uncertainty. Our primary setting is matching with transferable utilities, where the platform both matches agents and sets monetary transfers between them. We design an incentive-aware learning objective that captures the distance of a market outcome from equilibrium. Using this objective, we analyze the complexity of learning as a function of preference structure, casting learning as a stochastic multi-armed bandit problem. Algorithmically, we show that "optimism in the face of uncertainty," the principle underlying many bandit algorithms, applies to a primal-dual formulation of matching with transfers and leads to near-optimal regret bounds. Our work takes a first step toward elucidating when and how stable matchings arise in large, data-driven marketplaces.

[Towards Lower Bounds on the Depth of ReLU Neural Networks](#)

- Christoph Hertwich, Amitabh Basu, Marco Di Summa, Martin Skutella
- abstract@[open-review](#): We contribute to a better understanding of the class of functions that is represented by a neural network with ReLU activations and a given architecture. Using techniques from mixed-integer optimization, polyhedral theory, and tropical geometry, we provide a mathematical counterbalance to the universal approximation theorems which suggest that a single hidden layer is sufficient for learning tasks. In particular, we investigate whether the class of exactly representable functions strictly increases by adding more layers (with no restrictions on size). This problem has potential impact on algorithmic and statistical aspects because of the insight it provides into the class of functions represented by neural hypothesis classes. However, to the best of our knowledge, this question has not been investigated in the neural network literature. We also present upper bounds on the sizes of neural networks required to represent functions in these neural hypothesis classes.

[The Limitations of Large Width in Neural Networks: A Deep Gaussian Process Perspective](#)

- Geoff Pleiss, John P. Cunningham
- abstract@[open-review](#): Large width limits have been a recent focus of deep learning research: modulo computational practicalities, do wider networks outperform narrower ones? Answering this question has been challenging, as conventional networks gain representational power with width, potentially masking any negative effects. Our analysis in this paper decouples capacity and width via the generalization of neural networks to Deep Gaussian Processes (Deep GP), a class of nonparametric hierarchical models that subsume neural nets. In doing so, we aim to understand how width affects (standard) neural networks once they have sufficient capacity for a given modeling task. Our theoretical and empirical results on Deep GP suggest that large width can be detrimental to hierarchical models. Surprisingly, we prove that even nonparametric Deep GP converge to Gaussian processes, effectively becoming shallower without any increase in representational power. The posterior, which corresponds to a mixture of data-adaptable basis functions, becomes less data-dependent with width. Our tail analysis demonstrates that width and depth have opposite effects: depth accentuates a model's non-Gaussianity, while width makes models increasingly Gaussian. We find there is a "sweet spot" that maximizes test performance before the limiting GP behavior prevents adaptability, occurring at width = 1 or width = 2 for nonparametric Deep GP. These results make strong predictions about the same phenomenon in conventional neural networks trained with L2 regularization (analogous to a Gaussian prior on parameters): we show that such neural networks may need up to 500 – 1000 hidden units for sufficient capacity - depending on the dataset - but further width degrades performance.

[Exact marginal prior distributions of finite Bayesian neural networks](#)

- Jacob Zavatone-Veth, Cengiz Pehlevan
- abstract@[open-review](#): Bayesian neural networks are theoretically well-understood only in the infinite-width limit, where Gaussian priors over network weights yield Gaussian priors over network outputs. Recent work has suggested that finite Bayesian networks may outperform their infinite counterparts, but their non-Gaussian output priors have been characterized only through perturbative approaches. Here, we derive exact solutions for the function space priors for individual input examples of a class of finite fully-connected feedforward Bayesian neural networks. For deep linear networks, the prior has a simple expression in terms of the Meijer G-function. The prior of a finite ReLU network is a mixture of the priors of linear networks of smaller widths, corresponding to different numbers of active units in each layer. Our results unify previous descriptions of finite network priors in terms of their tail decay and large-width behavior.

[Spatiotemporal Joint Filter Decomposition in 3D Convolutional Neural Networks](#)

- Zichen Miao, Ze Wang, Xiuyuan Cheng, Qiang Qiu
- abstract@[open-review](#): In this paper, we introduce spatiotemporal joint filter decomposition to decouple spatial and temporal learning, while preserving spatiotemporal dependency in a video. A 3D convolutional filter is now jointly decomposed over a set of spatial and temporal filter atoms respectively. In this way, a 3D convolutional layer becomes three: a temporal atom layer, a spatial atom layer, and a joint coefficient layer, all three remaining convolutional. One obvious arithmetic manipulation allowed in our joint decomposition is to swap spatial or temporal atoms with a set of atoms that have the same number but different sizes, while keeping the remaining unchanged. For example, as shown later, we can now achieve tempo-invariance by simply dilating temporal atoms only. To illustrate this useful atom-swapping property, we further demonstrate how such a decomposition permits the direct learning of 3D CNNs with full-size videos through iterations of two consecutive sub-stages of learning: In the temporal stage, full-temporal downsampled-spatial data are used to learn temporal atoms and joint coefficients while fixing spatial atoms. In the spatial stage, full-spatial downsampled-temporal data are used for spatial atoms and joint coefficients while fixing temporal atoms. We show empirically on multiple action recognition datasets that, the decoupled spatiotemporal learning significantly reduces the model memory footprints, and allows deep 3D CNNs to model high-spatial long-temporal dependency with limited computational resources while delivering comparable performance.

[Pooling by Sliced-Wasserstein Embedding](#)

- Navid Naderizadeh, Joseph F Comer, Reed Andrews, Heiko Hoffmann, Soheil Kolouri
- abstract@[open-review](#): Learning representations from sets has become increasingly important with many applications in point cloud processing, graph learning, image/video recognition, and object detection. We introduce a geometrically-interpretable and generic pooling mechanism for aggregating a set of features into a fixed-dimensional representation. In particular, we treat elements of a set as samples from a probability distribution and propose an end-to-end trainable Euclidean embedding for sliced-Wasserstein distance to learn from set-structured data effectively. We evaluate our proposed pooling method on a wide variety of set-structured data, including point-cloud, graph, and image classification tasks, and demonstrate that our proposed method provides superior performance over existing set representation learning approaches. Our code is available at <https://github.com/navid-naderi/PSWE>.

[On the Theory of Reinforcement Learning with Once-per-Episode Feedback](#)

- Niladri Chatterji, Aldo Pacchiano, Peter Bartlett, Michael Jordan
- abstract@[open-review](#): We study a theory of reinforcement learning (RL) in which the learner receives binary feedback only once at the end of an episode. While this is an extreme test case for theory, it is also arguably more representative of real-world applications than the traditional requirement in RL practice that the learner receive feedback at every time step. Indeed, in many real-world applications of reinforcement learning, such as self-driving cars and robotics, it is easier to evaluate whether a learner's complete trajectory was either good" or bad," but harder to provide a reward signal at each step. To show that learning is possible in this more challenging setting, we study the case where trajectory labels are generated by an unknown parametric model, and provide a statistically and computationally efficient algorithm that achieves sublinear regret.

[ResNEsts and DenseNEsts: Block-based DNN Models with Improved Representation Guarantees](#)

- Kuan-Lin Chen, Ching-Hua Lee, Harinath Garudadri, Bhaskar D Rao
- abstract@[open-review](#): Models recently used in the literature proving residual networks (ResNets) are better than linear predictors are actually different from standard ResNets that have been widely used in computer vision. In addition to the assumptions such as scalar-valued output or single residual block, the models fundamentally considered in the literature have no nonlinearities at the final residual representation that feeds into the final affine layer. To codify such a difference in nonlinearities and reveal a linear estimation property, we define ResNEsts, i.e., Residual Nonlinear Estimators, by simply dropping nonlinearities at the last residual representation from standard ResNets. We show that wide ResNEsts with bottleneck blocks can always guarantee a very desirable training property that standard ResNets aim to achieve, i.e., adding more blocks does not decrease performance given the same set of basis elements. To prove that, we first recognize ResNEsts are basis function models that are limited by a coupling problem in basis learning and linear prediction. Then, to decouple prediction weights from basis learning, we construct a special architecture termed augmented ResNEst (A-ResNEst) that always guarantees no worse performance with the addition of a block. As a result, such an A-ResNEst establishes empirical risk lower bounds for a ResNEst using corresponding bases. Our results demonstrate ResNEsts indeed have a problem of diminishing feature reuse; however, it can be avoided by sufficiently expanding or widening the input space, leading to the above-mentioned desirable property. Inspired by the densely connected networks (DenseNets) that have been shown to outperform ResNets, we also propose a corresponding new model called Densely connected Nonlinear Estimator (DenseNEst). We show that any DenseNEst can be represented as a wide ResNEst with bottleneck blocks. Unlike ResNEsts, DenseNEsts exhibit the desirable property without any special architectural re-design.

[Locally private online change point detection](#)

- Tom Berrett, Yi Yu
- abstract@[open-review](#): We study online change point detection problems under the constraint of local differential privacy (LDP) where, in particular, the statistician does not have access to the raw data. As a concrete problem, we study a multivariate nonparametric regression problem. At each time point t , the raw data are assumed to be of the form (X_t, Y_t) , where X_t is a d -dimensional feature vector and Y_t is a response variable. Our primary aim is to detect changes in the regression function $m_t(x) = \mathbb{E}(Y_t | X_t=x)$ as soon as the change occurs. We provide algorithms which respect the LDP constraint, which control the false alarm probability, and which detect changes with a minimal (minimax rate-optimal) delay. To quantify the cost of privacy, we also present the optimal rate in the benchmark, non-private setting. These non-private results are also new to the literature and thus are interesting *per se*. In addition, we study the univariate mean online change point detection problem, under privacy constraints. This serves as the blueprint of studying more complicated private change point detection problems.

[Invariance Principle Meets Information Bottleneck for Out-of-Distribution Generalization](#)

- Kartik Ahuja, Ethan Caballero, Dinghuai Zhang, Jean-Christophe Gagnon-Audet, Yoshua Bengio, Ioannis Mitliagkas, Irina Rish
- abstract@[open-review](#): The invariance principle from causality is at the heart of notable approaches such as invariant risk minimization (IRM) that seek to address out-of-distribution (OOD) generalization failures. Despite the promising theory, invariance principle-based approaches fail in common classification tasks, where invariant (causal) features capture all the information about the label. Are these failures due to the methods failing to capture the invariance? Or is the invariance principle itself insufficient? To answer these questions, we revisit the fundamental assumptions in linear regression tasks, where invariance-based approaches were shown to provably generalize OOD. In contrast to the linear regression tasks, we show that for linear classification tasks we need much stronger restrictions on the distribution shifts, or otherwise OOD generalization is impossible. Furthermore, even with appropriate restrictions on distribution shifts in place, we show that the invariance principle alone is insufficient. We prove that a form of the information bottleneck constraint along with invariance helps address the key failures when invariant features capture all the information about the label and also retains the existing success when they do not. We propose an approach that incorporates both of these principles and demonstrate its effectiveness in several experiments.

[Repulsive Deep Ensembles are Bayesian](#)

- Francesco D'Angelo, Vincent Fortuin
- abstract@[open-review](#): Deep ensembles have recently gained popularity in the deep learning community for their conceptual simplicity and efficiency. However, maintaining functional diversity between ensemble members that are independently trained with gradient descent is challenging. This can lead to pathologies when adding more ensemble members, such as a saturation of the ensemble performance, which converges to the performance of a single model. Moreover, this does not only affect the quality of its predictions, but even more so the uncertainty estimates of the ensemble, and thus its performance on out-of-distribution data. We hypothesize that this limitation can be overcome by discouraging different ensemble members from collapsing to the same function. To this end, we introduce a kernelized repulsive term in the update rule of the deep ensembles. We show that this simple modification not only enforces and maintains diversity among the members but, even more importantly, transforms the maximum a posteriori inference into proper Bayesian inference. Namely, we show that the training dynamics of our proposed repulsive ensembles follow a Wasserstein gradient flow of the KL divergence with the true posterior. We study repulsive terms in weight and function space and empirically compare their performance to standard ensembles and Bayesian baselines on synthetic and real-world prediction tasks.

[BayesIMP: Uncertainty Quantification for Causal Data Fusion](#)

- Siu Lun Chau, Jean-Francois Ton, Javier González, Yee Teh, Dino Sejdinovic
- abstract@[open-review](#): While causal models are becoming one of the mainstays of machine learning, the problem of uncertainty quantification in causal inference remains challenging. In this paper, we study the causal data fusion problem, where data arising from multiple causal graphs are combined to estimate the average treatment effect of a target variable. As data arises from multiple sources and can vary in quality and sample size, principled uncertainty quantification becomes essential. To that end, we introduce *Bayesian Causal Mean Processes*, the framework which combines ideas from probabilistic integration and kernel mean embeddings to represent interventional distributions in the reproducing kernel Hilbert space, while taking into account the uncertainty within each causal graph. To demonstrate the informativeness of our uncertainty estimation, we apply our method to the Causal Bayesian Optimisation task and show improvements over state-of-the-art methods.

[RMM: Reinforced Memory Management for Class-Incremental Learning](#)

- Yaoyao Liu, Bernt Schiele, Qianru Sun
- abstract@[open-review](#): Class-Incremental Learning (CIL) [38] trains classifiers under a strict memory budget: in each incremental phase, learning is done for new data, most of which is abandoned to free space for the next phase. The preserved data are exemplars used for replaying. However, existing methods use a static and ad hoc strategy for memory allocation, which is often sub-optimal. In this work, we propose a dynamic memory management strategy that is optimized for the incremental phases and different object classes. We call our method reinforced memory management (RMM), leveraging reinforcement learning. RMM training is not naturally compatible with CIL as the past, and future data are strictly non-accessible during the incremental phases. We solve this by training the policy function of RMM on pseudo CIL tasks, e.g., the tasks built on the data of the zeroth phase, and then applying it to target tasks. RMM propagates two levels of actions: Level-1 determines how to split the memory between old and new classes, and Level-2 allocates memory for each specific class. In essence, it is an optimizable and general method for memory management that can be used in any replaying-based CIL method. For evaluation, we plug RMM into two top-performing baselines (LUCIR+AANets and POD+AANets [28]) and conduct experiments on three benchmarks (CIFAR-100, ImageNet-Subset, and ImageNet-Full). Our results show clear improvements, e.g., boosting POD+AANets by 3.6%, 4.4%, and 1.9% in the 25-Phase settings of the above benchmarks, respectively. The code is available at <https://class-il.mpi-inf.mpg.de/rmm/>.

[Learning Compact Representations of Neural Networks using Discriminative Masking \(DAM\)](#)

- Jie Bu, Arka Daw, M. Maruf, Anuj Karpatne
- abstract@[open-review](#): A central goal in deep learning is to learn compact representations of features at every layer of a neural network, which is useful for both unsupervised representation learning and structured network pruning. While there is a growing body of work in structured pruning, current state-of-the-art methods suffer from two key limitations: (i) instability during training, and (ii) need for an additional step of fine-tuning, which is resource-intensive. At the core of these limitations is the lack of a systematic approach that jointly prunes and refines weights during training in a single stage, and does not require any fine-tuning upon convergence to achieve state-of-the-art performance. We present a novel single-stage structured pruning method termed Discriminative Masking (DAM). The key intuition behind DAM is to discriminatively prefer some of the neurons to be refined during the training process, while gradually masking out other neurons. We show that our proposed DAM approach has remarkably good performance over a diverse range of applications in representation learning and structured pruning, including dimensionality reduction, recommendation system, graph representation learning, and structured pruning for image classification. We also theoretically show that the learning objective of DAM is directly related to minimizing the L_0 norm of the masking layer. All of our codes and datasets are available <https://github.com/jayroxis/dam-pytorch>.

[Neural Auto-Curricula in Two-Player Zero-Sum Games](#)

- Xidong Feng, Oliver Slumbers, Ziyu Wan, Bo Liu, Stephen McAleer, Ying Wen, Jun Wang, Yaodong Yang

- abstract@[open-review](#): When solving two-player zero-sum games, multi-agent reinforcement learning (MARL) algorithms often create populations of agents where, at each iteration, a new agent is discovered as the best response to a mixture over the opponent population. Within such a process, the update rules of "who to compete with" (i.e., the opponent mixture) and "how to beat them" (i.e., finding best responses) are underpinned by manually developed game theoretical principles such as fictitious play and Double Oracle. In this paper, we introduce a novel framework—Neural Auto-Curricula (NAC)—that leverages meta-gradient descent to automate the discovery of the learning update rule without explicit human design. Specifically, we parameterise the opponent selection module by neural networks and the best-response module by optimisation subroutines, and update their parameters solely via interaction with the game engine, where both players aim to minimise their exploitability. Surprisingly, even without human design, the discovered MARL algorithms achieve competitive or even better performance with the state-of-the-art population-based game solvers (e.g., PSRO) on Games of Skill, differentiable Lotto, non-transitive Mixture Games, Iterated Matching Pennies, and Kuhn Poker. Additionally, we show that NAC is able to generalise from small games to large games, for example training on Kuhn Poker and outperforming PSRO on Leduc Poker. Our work inspires a promising future direction to discover general MARL algorithms solely from data.

[ImageBART: Bidirectional Context with Multinomial Diffusion for Autoregressive Image Synthesis](#)

- Patrick Esser, Robin Rombach, Andreas Blattmann, Bjorn Ommer
- abstract@[open-review](#): Autoregressive models and their sequential factorization of the data likelihood have recently demonstrated great potential for image representation and synthesis. Nevertheless, they incorporate image context in a linear 1D order by attending only to previously synthesized image patches above or to the left. Not only is this unidirectional, sequential bias of attention unnatural for images as it disregards large parts of a scene until synthesis is almost complete. It also processes the entire image on a single scale, thus ignoring more global contextual information up to the gist of the entire scene. As a remedy we incorporate a coarse-to-fine hierarchy of context by combining the autoregressive formulation with a multinomial diffusion process: Whereas a multistage diffusion process successively compresses and removes information to coarsen an image, we train a Markov chain to invert this process. In each stage, the resulting autoregressive ImageBART model progressively incorporates context from previous stages in a coarse-to-fine manner. Experiments demonstrate the gain over current autoregressive models, continuous diffusion probabilistic models, and latent variable models. Moreover, the approach enables to control the synthesis process and to trade compression rate against reconstruction accuracy, while still guaranteeing visually plausible results.

[From global to local MDI variable importances for random forests and when they are Shapley values](#)

- Antonio Sutera, Gilles Louppe, Van Anh Huynh-Thu, Louis Wehenkel, Pierre Geurts
- abstract@[open-review](#): Random forests have been widely used for their ability to provide so-called importance measures, which give insight at a global (per dataset) level on the relevance of input variables to predict a certain output. On the other hand, methods based on Shapley values have been introduced to refine the analysis of feature relevance in tree-based models to a local (per instance) level. In this context, we first show that the global Mean Decrease of Impurity (MDI) variable importance scores correspond to Shapley values under some conditions. Then, we derive a local MDI importance measure of variable relevance, which has a very natural connection with the global MDI measure and can be related to a new notion of local feature relevance. We further link local MDI importances with Shapley values and discuss them in the light of related measures from the literature. The measures are illustrated through experiments on several classification and regression problems.

[Adversarial Robustness of Streaming Algorithms through Importance Sampling](#)

- Vladimir Braverman, Avinatan Hassidim, Yossi Matias, Mariano Schain, Sandeep Silwal, Samson Zhou
- abstract@[open-review](#): Robustness against adversarial attacks has recently been at the forefront of algorithmic design for machine learning tasks. In the adversarial streaming model, an adversary gives an algorithm a sequence of adaptively chosen updates u_1, \dots, u_n as a data stream. The goal of the algorithm is to compute or approximate some predetermined function for every prefix of the adversarial stream, but the adversary may generate future updates based on previous outputs of the algorithm. In particular, the adversary may gradually learn the random bits internally used by an algorithm to manipulate dependencies in the input. This is especially problematic as many important problems in the streaming model require randomized algorithms, as they are known to not admit any deterministic algorithms that use sublinear space. In this paper, we introduce adversarially robust streaming algorithms for central machine learning and algorithmic tasks, such as regression and clustering, as well as their more general counterparts, subspace embedding, low-rank approximation, and coresnet construction. For regression and other numerical linear algebra related tasks, we consider the row arrival streaming model. Our results are based on a simple, but powerful, observation that many importance sampling-based algorithms give rise to adversarial robustness which is in contrast to sketching based algorithms, which are very prevalent in the streaming literature but suffer from adversarial attacks. In addition, we show that the well-known merge and reduce paradigm in streaming is adversarially robust. Since the merge and reduce paradigm allows coresnet constructions in the streaming setting, we thus obtain robust algorithms for k -means, k -median, k -center, Bregman clustering, projective clustering, principal component analysis (PCA) and non-negative matrix factorization. To the best of our knowledge, these are the first adversarially robust results for these problems yet require no new algorithmic implementations. Finally, we empirically confirm the robustness of our algorithms on various adversarial attacks and demonstrate that by contrast, some common existing algorithms are not robust.

[Tractable Regularization of Probabilistic Circuits](#)

- Anji Liu, Guy Van den Broeck
- abstract@[open-review](#): Probabilistic Circuits (PCs) are a promising avenue for probabilistic modeling. They combine advantages of probabilistic graphical models (PGMs) with those of neural networks (NNs). Crucially, however, they are tractable probabilistic models, supporting efficient and exact computation of many probabilistic inference queries, such as marginals and MAP. Further, since PCs are structured computation graphs, they can take advantage of deep-learning-style parameter updates, which greatly improves their scalability. However, this innovation also makes PCs prone to overfitting, which has been observed in many standard benchmarks. Despite the existence of abundant regularization techniques for both PGMs and NNs, they are not effective enough when applied to PCs. Instead, we re-think regularization for PCs and propose two intuitive techniques, data softening and entropy regularization, that both take advantage of PCs' tractability and still have an efficient implementation as a computation graph. Specifically, data softening provides a principled way to add uncertainty in datasets in closed form, which implicitly regularizes PC parameters. To learn parameters from a softened dataset, PCs only need linear time by virtue of their tractability. In entropy regularization, the exact entropy of the distribution encoded by a PC can be regularized directly, which is again infeasible for most other density estimation models. We show that both methods consistently improve the generalization performance of a wide variety of PCs. Moreover, when paired with a simple PC structure, we achieved state-of-the-art results on 10 out of 20 standard discrete density estimation benchmarks. Open-source code and experiments are available at <https://github.com/UCLA-StarAI/Tractable-PC-Regularization>.

[On Interaction Between Augmentations and Corruptions in Natural Corruption Robustness](#)

- Eric Mintun, Alexander Kirillov, Saining Xie
- abstract@[open-review](#): Invariance to a broad array of image corruptions, such as warping, noise, or color shifts, is an important aspect of building robust models in computer vision. Recently, several new data augmentations have been proposed that significantly improve performance on ImageNet-C, a benchmark of such corruptions. However, there is still a lack of basic understanding on the relationship between data augmentations and test-time corruptions. To this end, we develop a feature space for image transforms, and then use a new measure in this space between augmentations and corruptions called the Minimal Sample Distance to demonstrate there is a strong correlation between similarity and performance. We then investigate recent data augmentations and observe a significant degradation in corruption robustness when the test-time corruptions are sampled to be perceptually

dissimilar from ImageNet-C in this feature space. Our results suggest that test error can be improved by training on perceptually similar augmentations, and data augmentations may not generalize well beyond the existing benchmark. We hope our results and tools will allow for more robust progress towards improving robustness to image corruptions. We provide code at <https://github.com/facebookresearch/augmentation-corruption>.

[Dynamic Distillation Network for Cross-Domain Few-Shot Recognition with Unlabeled Data](#)

- Ashraful Islam, Chun-Fu (Richard) Chen, Rameswar Panda, Leonid Karlinsky, Rogerio Feris, Richard J. Radke
- abstract@[open-review](#): Most existing works in few-shot learning rely on meta-learning the network on a large base dataset which is typically from the same domain as the target dataset. We tackle the problem of cross-domain few-shot learning where there is a large shift between the base and target domain. The problem of cross-domain few-shot recognition with unlabeled target data is largely unaddressed in the literature. STARTUP was the first method that tackles this problem using self-training. However, it uses a fixed teacher pretrained on a labeled base dataset to create soft labels for the unlabeled target samples. As the base dataset and unlabeled dataset are from different domains, projecting the target images in the class-domain of the base dataset with a fixed pretrained model might be sub-optimal. We propose a simple dynamic distillation-based approach to facilitate unlabeled images from the novel/base dataset. We impose consistency regularization by calculating predictions from the weakly-augmented versions of the unlabeled images from a teacher network and matching it with the strongly augmented versions of the same images from a student network. The parameters of the teacher network are updated as exponential moving average of the parameters of the student network. We show that the proposed network learns representation that can be easily adapted to the target domain even though it has not been trained with target-specific classes during the pretraining phase. Our model outperforms the current state-of-the art method by 4.4% for 1-shot and 3.6% for 5-shot classification in the BSCD-FSL benchmark, and also shows competitive performance on traditional in-domain few-shot learning task.

[Hypergraph Propagation and Community Selection for Objects Retrieval](#)

- Guoyuan An, Yuchi Huo, Sung-eui Yoon
- abstract@[open-review](#): Spatial verification is a crucial technique for particular object retrieval. It utilizes spatial information for the accurate detection of true positive images. However, existing query expansion and diffusion methods cannot efficiently propagate the spatial information in an ordinary graph with scalar edge weights, resulting in low recall or precision. To tackle these problems, we propose a novel hypergraph-based framework that efficiently propagates spatial information in query time and retrieves an object in the database accurately. Additionally, we propose using the image graph's structure information through community selection technique, to measure the accuracy of the initial search result and to provide correct starting points for hypergraph propagation without heavy spatial verification computations. Experiment results on ROxford and RParis show that our method significantly outperforms the existing query expansion and diffusion methods.

[Deep learning is adaptive to intrinsic dimensionality of model smoothness in anisotropic Besov space](#)

- Taiji Suzuki, Atsushi Nitanda
- abstract@[open-review](#): Deep learning has exhibited superior performance for various tasks, especially for high-dimensional datasets, such as images. To understand this property, we investigate the approximation and estimation ability of deep learning on {\it anisotropic Besov spaces}. The anisotropic Besov space is characterized by direction-dependent smoothness and includes several function classes that have been investigated thus far. We demonstrate that the approximation error and estimation error of deep learning only depend on the average value of the smoothness parameters in all directions. Consequently, the curse of dimensionality can be avoided if the smoothness of the target function is highly anisotropic. Unlike existing studies, our analysis does not require a low-dimensional structure of the input data. We also investigate the minimax optimality of deep learning and compare its performance with that of the kernel method (more generally, linear estimators). The results show that deep learning has better dependence on the input dimensionality if the target function possesses anisotropic smoothness, and it achieves an adaptive rate for functions with spatially inhomogeneous smoothness.

[QuPeD: Quantized Personalization via Distillation with Applications to Federated Learning](#)

- Kaan Ozkara, Navjot Singh, Deepesh Data, Suhas Diggavi
- abstract@[open-review](#): Traditionally, federated learning (FL) aims to train a single global model while collaboratively using multiple clients and a server. Two natural challenges that FL algorithms face are heterogeneity in data across clients and collaboration of clients with diverse resources. In this work, we introduce a quantized and personalized FL algorithm QuPeD that facilitates collective (personalized model compression) training via knowledge distillation (KD) among clients who have access to heterogeneous data and resources. For personalization, we allow clients to learn compressed personalized models with different quantization parameters and model dimensions/structures. Towards this, first we propose an algorithm for learning quantized models through a relaxed optimization problem, where quantization values are also optimized over. When each client participating in the (federated) learning process has different requirements of the compressed model (both in model dimension and precision), we formulate a compressed personalization framework by introducing knowledge distillation loss for local client objectives collaborating through a global model. We develop an alternating proximal gradient update for solving this compressed personalization problem, and analyze its convergence properties. Numerically, we validate that QuPeD outperforms competing personalized FL methods, FedAvg, and local training of clients in various heterogeneous settings.

[Model Adaptation: Historical Contrastive Learning for Unsupervised Domain Adaptation without Source Data](#)

- Jiaxing Huang, Dayan Guan, Aoran Xiao, Shijian Lu
- abstract@[open-review](#): Unsupervised domain adaptation aims to align a labeled source domain and an unlabeled target domain, but it requires to access the source data which often raises concerns in data privacy, data portability and data transmission efficiency. We study unsupervised model adaptation (UMA), or called Unsupervised Domain Adaptation without Source Data, an alternative setting that aims to adapt source-trained models towards target distributions without accessing source data. To this end, we design an innovative historical contrastive learning (HCL) technique that exploits historical source hypothesis to make up for the absence of source data in UMA. HCL addresses the UMA challenge from two perspectives. First, it introduces historical contrastive instance discrimination (HCID) that learns from target samples by contrasting their embeddings which are generated by the currently adapted model and the historical models. With the historical models, HCID encourages UMA to learn instance-discriminative target representations while preserving the source hypothesis. Second, it introduces historical contrastive category discrimination (HCCD) that pseudo-labels target samples to learn category-discriminative target representations. Specifically, HCCD re-weights pseudo labels according to their prediction consistency across the current and historical models. Extensive experiments show that HCL outperforms and state-of-the-art methods consistently across a variety of visual tasks and setups.

[The Out-of-Distribution Problem in Explainability and Search Methods for Feature Importance Explanations](#)

- Peter Hase, Harry Xie, Mohit Bansal
- abstract@[open-review](#): Feature importance (FI) estimates are a popular form of explanation, and they are commonly created and evaluated by computing the change in model confidence caused by removing certain input features at test time. For example, in the standard Sufficiency metric, only the top-k most important tokens are kept. In this paper, we study several under-explored dimensions of FI explanations, providing conceptual and empirical improvements for this form of explanation. First, we advance a new argument for why it can be problematic to remove features from an input when creating or evaluating explanations: the fact that these counterfactual inputs are out-of-distribution (OOD) to models implies that the resulting explanations are socially misaligned. The crux of the problem is that the model prior and random weight initialization influence the explanations (and

explanation metrics) in unintended ways. To resolve this issue, we propose a simple alteration to the model training process, which results in more socially aligned explanations and metrics. Second, we compare among five approaches for removing features from model inputs. We find that some methods produce more OOD counterfactuals than others, and we make recommendations for selecting a feature-replacement function. Finally, we introduce four search-based methods for identifying FI explanations and compare them to strong baselines, including LIME, Anchors, and Integrated Gradients. Through experiments with six diverse text classification datasets, we find that the only method that consistently outperforms random search is a Parallel Local Search (PLS) that we introduce. Improvements over the second best method are as large as 5.4 points for Sufficiency and 17 points for Comprehensiveness.

[Control Variates for Slate Off-Policy Evaluation](#)

- Nikos Vlassis, Ashok Chandrashekhar, Fernando Amat, Nathan Kallus
- abstract@[open-review](#): We study the problem of off-policy evaluation from batched contextual bandit data with multidimensional actions, often termed slates. The problem is common to recommender systems and user-interface optimization, and it is particularly challenging because of the combinatorially-sized action space. Swaminathan et al. (2017) have proposed the pseudoinverse (PI) estimator under the assumption that the conditional mean rewards are additive in actions. Using control variates, we consider a large class of unbiased estimators that includes as specific cases the PI estimator and (asymptotically) its self-normalized variant. By optimizing over this class, we obtain new estimators with risk improvement guarantees over both the PI and the self-normalized PI estimators. Experiments with real-world recommender data as well as synthetic data validate these improvements in practice.

[Stabilizing Deep Q-Learning with ConvNets and Vision Transformers under Data Augmentation](#)

- Nicklas Hansen, Hao Su, Xiaolong Wang
- abstract@[open-review](#): While agents trained by Reinforcement Learning (RL) can solve increasingly challenging tasks directly from visual observations, generalizing learned skills to novel environments remains very challenging. Extensive use of data augmentation is a promising technique for improving generalization in RL, but it is often found to decrease sample efficiency and can even lead to divergence. In this paper, we investigate causes of instability when using data augmentation in common off-policy RL algorithms. We identify two problems, both rooted in high-variance Q-targets. Based on our findings, we propose a simple yet effective technique for stabilizing this class of algorithms under augmentation. We perform extensive empirical evaluation of image-based RL using both ConvNets and Vision Transformers (ViT) on a family of benchmarks based on DeepMind Control Suite, as well as in robotic manipulation tasks. Our method greatly improves stability and sample efficiency of ConvNets under augmentation, and achieves generalization results competitive with state-of-the-art methods for image-based RL in environments with unseen visuals. We further show that our method scales to RL with ViT-based architectures, and that data augmentation may be especially important in this setting.

[On Effective Scheduling of Model-based Reinforcement Learning](#)

- Hang Lai, Jian Shen, Weinan Zhang, Yimin Huang, Xing Zhang, Ruiming Tang, Yong Yu, Zhenguo Li
- abstract@[open-review](#): Model-based reinforcement learning has attracted wide attention due to its superior sample efficiency. Despite its impressive success so far, it is still unclear how to appropriately schedule the important hyperparameters to achieve adequate performance, such as the real data ratio for policy optimization in Dyna-style model-based algorithms. In this paper, we first theoretically analyze the role of real data in policy training, which suggests that gradually increasing the ratio of real data yields better performance. Inspired by the analysis, we propose a framework named AutoMBPO to automatically schedule the real data ratio as well as other hyperparameters in training model-based policy optimization (MBPO) algorithm, a representative running case of model-based methods. On several continuous control tasks, the MBPO instance trained with hyperparameters scheduled by AutoMBPO can significantly surpass the original one, and the real data ratio schedule found by AutoMBPO shows consistency with our theoretical analysis.

[Removing Inter-Experimental Variability from Functional Data in Systems Neuroscience](#)

- Dominic Gonschorek, Larissa Höfling, Klaudia P. Szatkó, Katrin Franke, Timm Schubert, Benjamin Dunn, Philipp Berens, David Klindt, Thomas Euler
- abstract@[open-review](#): Integrating data from multiple experiments is common practice in systems neuroscience but it requires inter-experimental variability to be negligible compared to the biological signal of interest. This requirement is rarely fulfilled; systematic changes between experiments can drastically affect the outcome of complex analysis pipelines. Modern machine learning approaches designed to adapt models across multiple data domains offer flexible ways of removing inter-experimental variability where classical statistical methods often fail. While applications of these methods have been mostly limited to single-cell genomics, in this work, we develop a theoretical framework for domain adaptation in systems neuroscience. We implement this in an adversarial optimization scheme that removes inter-experimental variability while preserving the biological signal. We compare our method to previous approaches on a large-scale dataset of two-photon imaging recordings of retinal bipolar cell responses to visual stimuli. This dataset provides a unique benchmark as it contains biological signal from well-defined cell types that is obscured by large inter-experimental variability. In a supervised setting, we compare the generalization performance of cell type classifiers across experiments, which we validate with anatomical cell type distributions from electron microscopy data. In an unsupervised setting, we remove inter-experimental variability from the data which can then be fed into arbitrary downstream analyses. In both settings, we find that our method achieves the best trade-off between removing inter-experimental variability and preserving biological signal. Thus, we offer a flexible approach to remove inter-experimental variability and integrate datasets across experiments in systems neuroscience. Code available at <https://github.com/eulerlab/rave>.

[Learning Knowledge Graph-based World Models of Textual Environments](#)

- Prithviraj Ammanabrolu, Mark Riedl
- abstract@[open-review](#): World models improve a learning agent's ability to efficiently operate in interactive and situated environments. This work focuses on the task of building world models of text-based game environments. Text-based games, or interactive narratives, are reinforcement learning environments in which agents perceive and interact with the world using textual natural language. These environments contain long, multi-step puzzles or quests woven through a world that is filled with hundreds of characters, locations, and objects. Our world model learns to simultaneously: (1) predict changes in the world caused by an agent's actions when representing the world as a knowledge graph; and (2) generate the set of contextually relevant natural language actions required to operate in the world. We frame this task as a Set of Sequences generation problem by exploiting the inherent structure of knowledge graphs and actions and introduce both a transformer-based multi-task architecture and a loss function to train it. A zero-shot ablation study on never-before-seen textual worlds shows that our methodology significantly outperforms existing textual world modeling techniques as well as the importance of each of our contributions.

[Damped Anderson Mixing for Deep Reinforcement Learning: Acceleration, Convergence, and Stabilization](#)

- Ke Sun, Yafei Wang, Yi Liu, yingnan zhao, Bo Pan, Shangling Jui, Bei Jiang, Linglong Kong
- abstract@[open-review](#): Anderson mixing has been heuristically applied to reinforcement learning (RL) algorithms for accelerating convergence and improving the sampling efficiency of deep RL. Despite its heuristic improvement of convergence, a rigorous mathematical justification for the benefits of Anderson mixing in RL has not yet been put forward. In this paper, we provide deeper insights into a class of acceleration schemes built on Anderson mixing that improve the convergence of deep RL algorithms. Our main results establish a connection between Anderson mixing and quasi-Newton methods and prove that Anderson mixing increases the convergence radius of policy iteration schemes by an extra contraction factor. The key focus of the analysis roots in the fixed-point iteration nature of RL. We further propose a stabilization strategy by introducing a stable regularization term in Anderson

mixing and a differentiable, non-expansive MellowMax operator that can allow both faster convergence and more stable behavior. Extensive experiments demonstrate that our proposed method enhances the convergence, stability, and performance of RL algorithms.

[Approximate Decomposable Submodular Function Minimization for Cardinality-Based Components](#)

- Nate Veldt, Austin R. Benson, Jon Kleinberg
- abstract@[open-review](#): Minimizing a sum of simple submodular functions of limited support is a special case of general submodular function minimization that has seen numerous applications in machine learning. We develop faster techniques for instances where components in the sum are cardinality-based, meaning they depend only on the size of the input set. This variant is one of the most widely applied in practice, encompassing, e.g., common energy functions arising in image segmentation and recent generalized hypergraph cut functions. We develop the first approximation algorithms for this problem, where the approximations can be quickly computed via reduction to a sparse graph cut problem, with graph sparsity controlled by the desired approximation factor. Our method relies on a new connection between sparse graph reduction techniques and piecewise linear approximations to concave functions. Our sparse reduction technique leads to significant improvements in theoretical runtimes, as well as substantial practical gains in problems ranging from benchmark image segmentation tasks to hypergraph clustering problems.

[Episodic Multi-agent Reinforcement Learning with Curiosity-driven Exploration](#)

- Lulu Zheng, Jiarui Chen, Jianhao Wang, Jiamin He, Yujing Hu, Yingfeng Chen, Changjie Fan, Yang Gao, Chongjie Zhang
- abstract@[open-review](#): Efficient exploration in deep cooperative multi-agent reinforcement learning (MARL) still remains challenging in complex coordination problems. In this paper, we introduce a novel Episodic Multi-agent reinforcement learning with Curiosity-driven exploration, called EMC. We leverage an insight of popular factorized MARL algorithms that the ``induced'' individual Q-values, i.e., the individual utility functions used for local execution, are the embeddings of local action-observation histories, and can capture the interaction between agents due to reward backpropagation during centralized training. Therefore, we use prediction errors of individual Q-values as intrinsic rewards for coordinated exploration and utilize episodic memory to exploit explored informative experience to boost policy training. As the dynamics of an agent's individual Q-value function captures the novelty of states and the influence from other agents, our intrinsic reward can induce coordinated exploration to new or promising states. We illustrate the advantages of our method by didactic examples, and demonstrate its significant outperformance over state-of-the-art MARL baselines on challenging tasks in the StarCraft II micromanagement benchmark.

[Two Sides of Meta-Learning Evaluation: In vs. Out of Distribution](#)

- Amrith Setlur, Oscar Li, Virginia Smith
- abstract@[open-review](#): We categorize meta-learning evaluation into two settings: \$\textit{in-distribution}\$ [ID], in which the train and test tasks are sampled \$\textit{iid}\$ from the same underlying task distribution, and \$\textit{out-of-distribution}\$ [OOD], in which they are not. While most meta-learning theory and some FSL applications follow the ID setting, we identify that most existing few-shot classification benchmarks instead reflect OOD evaluation, as they use disjoint sets of train (base) and test (novel) classes for task generation. This discrepancy is problematic because -- as we show on numerous benchmarks -- meta-learning methods that perform better on existing OOD datasets may perform significantly worse in the ID setting. In addition, in the OOD setting, even though current FSL benchmarks seem befitting, our study highlights concerns in 1) reliably performing model selection for a given meta-learning method, and 2) consistently comparing the performance of different methods. To address these concerns, we provide suggestions on how to construct FSL benchmarks to allow for ID evaluation as well as more reliable OOD evaluation. Our work aims to inform the meta-learning community about the importance and distinction of ID vs. OOD evaluation, as well as the subtleties of OOD evaluation with current benchmarks.

[Debiased Visual Question Answering from Feature and Sample Perspectives](#)

- Zhiqian Wen, Guanghui Xu, Mingkui Tan, Qingyao Wu, Qi Wu
- abstract@[open-review](#): Visual question answering (VQA) is designed to examine the visual-textual reasoning ability of an intelligent agent. However, recent observations show that many VQA models may only capture the biases between questions and answers in a dataset rather than showing real reasoning abilities. For example, given a question, some VQA models tend to output the answer that occurs frequently in the dataset and ignore the images. To reduce this tendency, existing methods focus on weakening the language bias. Meanwhile, only a few works also consider vision bias implicitly. However, these methods introduce additional annotations or show unsatisfactory performance. Moreover, not all biases are harmful to the models. Some “biases” learnt from datasets represent natural rules of the world and can help limit the range of answers. Thus, how to filter and remove the true negative biases in language and vision modalities remain a major challenge. In this paper, we propose a method named D-VQA to alleviate the above challenges from the feature and sample perspectives. Specifically, from the feature perspective, we build a question-to-answer and vision-to-answer branch to capture the language and vision biases, respectively. Next, we apply two unimodal bias detection modules to explicitly recognise and remove the negative biases. From the sample perspective, we construct two types of negative samples to assist the training of the models, without introducing additional annotations. Extensive experiments on the VQA-CP v2 and VQA v2 datasets demonstrate the effectiveness of our D-VQA method.

[Towards a Unified Game-Theoretic View of Adversarial Perturbations and Robustness](#)

- Jie Ren, Die Zhang, Yisen Wang, Lu Chen, Zhanpeng Zhou, Yiting Chen, Xu Cheng, Xin Wang, Meng Zhou, Jie Shi, Quanshi Zhang
- abstract@[open-review](#): This paper provides a unified view to explain different adversarial attacks and defense methods, i.e. the view of multi-order interactions between input variables of DNNs. Based on the multi-order interaction, we discover that adversarial attacks mainly affect high-order interactions to fool the DNN. Furthermore, we find that the robustness of adversarially trained DNNs comes from category-specific low-order interactions. Our findings provide a potential method to unify adversarial perturbations and robustness, which can explain the existing robustness-boosting methods in a principle way. Besides, our findings also make a revision of previous inaccurate understanding of the shape bias of adversarially learned features. Our code is available online at <https://github.com/Jie-Ren/A-Unified-Game-Theoretic-Interpretation-of-Adversarial-Robustness>.

[On the Out-of-distribution Generalization of Probabilistic Image Modelling](#)

- Mingtian Zhang, Andi Zhang, Steven McDonagh
- abstract@[open-review](#): Out-of-distribution (OOD) detection and lossless compression constitute two problems that can be solved by the training of probabilistic models on a first dataset with subsequent likelihood evaluation on a second dataset, where data distributions differ. By defining the generalization of probabilistic models in terms of likelihood we show that, in the case of image models, the OOD generalization ability is dominated by local features. This motivates our proposal of a Local Autoregressive model that exclusively models local image features towards improving OOD performance. We apply the proposed model to OOD detection tasks and achieve state-of-the-art unsupervised OOD detection performance without the introduction of additional data. Additionally, we employ our model to build a new lossless image compressor: NeLLoC (Neural Local Lossless Compressor) and report state-of-the-art compression rates and model size.

[Exploiting Local Convergence of Quasi-Newton Methods Globally: Adaptive Sample Size Approach](#)

- Qiujiang Jin, Aryan Mokhtari

- abstract@[open-review](#): In this paper, we study the application of quasi-Newton methods for solving empirical risk minimization (ERM) problems defined over a large dataset. Traditional deterministic and stochastic quasi-Newton methods can be executed to solve such problems; however, it is known that their global convergence rate may not be better than first-order methods, and their local superlinear convergence only appears towards the end of the learning process. In this paper, we use an adaptive sample size scheme that exploits the superlinear convergence of quasi-Newton methods globally and throughout the entire learning process. The main idea of the proposed adaptive sample size algorithms is to start with a small subset of data points and solve their corresponding ERM problem within its statistical accuracy, and then enlarge the sample size geometrically and use the optimal solution of the problem corresponding to the smaller set as an initial point for solving the subsequent ERM problem with more samples. We show that if the initial sample size is sufficiently large and we use quasi-Newton methods to solve each subproblem, the subproblems can be solved superlinearly fast (after at most three iterations), as we guarantee that the iterates always stay within a neighborhood that quasi-Newton methods converge superlinearly. Numerical experiments on various datasets confirm our theoretical results and demonstrate the computational advantages of our method.

[**PDE-GCN: Novel Architectures for Graph Neural Networks Motivated by Partial Differential Equations**](#)

- Moshe Eliasof, Eldad Haber, Eran Treister
- abstract@[open-review](#): Graph neural networks are increasingly becoming the go-to approach in various fields such as computer vision, computational biology and chemistry, where data are naturally explained by graphs. However, unlike traditional convolutional neural networks, deep graph networks do not necessarily yield better performance than shallow graph networks. This behavior usually stems from the over-smoothing phenomenon. In this work, we propose a family of architectures to control this behavior by design. Our networks are motivated by numerical methods for solving Partial Differential Equations (PDEs) on manifolds, and as such, their behavior can be explained by similar analysis. Moreover, as we demonstrate using an extensive set of experiments, our PDE-motivated networks can generalize and be effective for various types of problems from different fields. Our architectures obtain better or on par with the current state-of-the-art results for problems that are typically approached using different architectures.

[**Information Directed Reward Learning for Reinforcement Learning**](#)

- David Lindner, Matteo Turchetta, Sebastian Tschiatschek, Kamil Ciosek, Andreas Krause
- abstract@[open-review](#): For many reinforcement learning (RL) applications, specifying a reward is difficult. In this paper, we consider an RL setting where the agent can obtain information about the reward only by querying an expert that can, for example, evaluate individual states or provide binary preferences over trajectories. From such expensive feedback, we aim to learn a model of the reward function that allows standard RL algorithms to achieve high expected return with as few expert queries as possible. For this purpose, we propose Information Directed Reward Learning (IDRL), which uses a Bayesian model of the reward function and selects queries that maximize the information gain about the difference in return between potentially optimal policies. In contrast to prior active reward learning methods designed for specific types of queries, IDRL naturally accommodates different query types. Moreover, by shifting the focus from reducing the reward approximation error to improving the policy induced by the reward model, it achieves similar or better performance with significantly fewer queries. We support our findings with extensive evaluations in multiple environments and with different types of queries.

[**SSMF: Shifting Seasonal Matrix Factorization**](#)

- Koki Kawabata, Siddharth Bhatia, Rui Liu, Mohit Wadhwa, Bryan Hooi
- abstract@[open-review](#): Given taxi-ride counts information between departure and destination locations, how can we forecast their future demands? In general, given a data stream of events with seasonal patterns that innovate over time, how can we effectively and efficiently forecast future events? In this paper, we propose Shifting Seasonal Matrix Factorization approach, namely SSMF, that can adaptively learn multiple seasonal patterns (called regimes), as well as switching between them. Our proposed method has the following properties: (a) it accurately forecasts future events by detecting regime shifts in seasonal patterns as the data stream evolves; (b) it works in an online setting, i.e., processes each observation in constant time and memory; (c) it effectively realizes regime shifts without human intervention by using a lossless data compression scheme. We demonstrate that our algorithm outperforms state-of-the-art baseline methods by accurately forecasting upcoming events on three real-world data streams.

[**Associative Memories via Predictive Coding**](#)

- Tommaso Salvatori, Yuhang Song, Yujian Hong, Lei Sha, Simon Frieder, Zhenghua Xu, Rafal Bogacz, Thomas Lukasiewicz
- abstract@[open-review](#): Associative memories in the brain receive and store patterns of activity registered by the sensory neurons, and are able to retrieve them when necessary. Due to their importance in human intelligence, computational models of associative memories have been developed for several decades now. In this paper, we present a novel neural model for realizing associative memories, which is based on a hierarchical generative network that receives external stimuli via sensory neurons. It is trained using predictive coding, an error-based learning algorithm inspired by information processing in the cortex. To test the model's capabilities, we perform multiple retrieval experiments from both corrupted and incomplete data points. In an extensive comparison, we show that this new model outperforms in retrieval accuracy and robustness popular associative memory models, such as autoencoders trained via backpropagation, and modern Hopfield networks. In particular, in completing partial data points, our model achieves remarkable results on natural image datasets, such as ImageNet, with a surprisingly high accuracy, even when only a tiny fraction of pixels of the original images is presented. Our model provides a plausible framework to study learning and retrieval of memories in the brain, as it closely mimics the behavior of the hippocampus as a memory index and generative model.

[**Robust and differentially private mean estimation**](#)

- Xiyang Liu, Weihao Kong, Sham Kakade, Sewoong Oh
- abstract@[open-review](#): In statistical learning and analysis from shared data, which is increasingly widely adopted in platforms such as federated learning and meta-learning, there are two major concerns: privacy and robustness. Each participating individual should be able to contribute without the fear of leaking one's sensitive information. At the same time, the system should be robust in the presence of malicious participants inserting corrupted data. Recent algorithmic advances in learning from shared data focus on either one of these threats, leaving the system vulnerable to the other. We bridge this gap for the canonical problem of estimating the mean from i.i.d.~samples. We introduce PRIME, which is the first efficient algorithm that achieves both privacy and robustness for a wide range of distributions. We further complement this result with a novel exponential time algorithm that improves the sample complexity of PRIME, achieving a near-optimal guarantee and matching that of a known lower bound for (non-robust) private mean estimation. This proves that there is no extra statistical cost to simultaneously guaranteeing privacy and robustness.

[**Adaptable Agent Populations via a Generative Model of Policies**](#)

- Kenneth Derek, Phillip Isola
- abstract@[open-review](#): In the natural world, life has found innumerable ways to survive and often thrive. Between and even within species, each individual is in some manner unique, and this diversity lends adaptability and robustness to life. In this work, we aim to learn a space of diverse and high-reward policies in a given environment. To this end, we introduce a generative model of policies for reinforcement learning, which maps a low-dimensional latent space to an agent policy space. Our method enables learning an entire population of agent policies, without requiring the use of separate policy parameters. Just as real world populations can adapt and evolve via natural selection, our method is able to adapt to changes in our environment solely by selecting for policies in latent space. We test our generative model's capabilities in a variety of environments, including an open-

ended grid-world and a two-player soccer environment. Code, visualizations, and additional experiments can be found at <https://kennyderek.github.io/adap/>.

[A No-go Theorem for Robust Acceleration in the Hyperbolic Plane](#)

- Linus Hamilton, Ankur Moitra
- abstract@[open-review](#): In recent years there has been significant effort to adapt the key tools and ideas in convex optimization to the Riemannian setting. One key challenge has remained: Is there a Nesterov-like accelerated gradient method for geodesically convex functions on a Riemannian manifold? Recent work has given partial answers and the hope was that this ought to be possible. Here we prove that in a noisy setting, there is no analogue of accelerated gradient descent for geodesically convex functions on the hyperbolic plane. Our results apply even when the noise is exponentially small. The key intuition behind our proof is short and simple: In negatively curved spaces, the volume of a ball grows so fast that information about the past gradients is not useful in the future.

[Privately Learning Mixtures of Axis-Aligned Gaussians](#)

- Ishaq Aden-Ali, Hassan Ashtiani, Christopher Liaw
- abstract@[open-review](#): We consider the problem of learning multivariate Gaussians under the constraint of approximate differential privacy. We prove that $\widetilde{O}(k^2 d \log^{3/2}(1/\delta) / \alpha^2 \epsilon)$ samples are sufficient to learn a mixture of k axis-aligned Gaussians in \mathbb{R}^d to within total variation distance α while satisfying (ϵ, δ) -differential privacy. This is the first result for privately learning mixtures of unbounded axis-aligned (or even unbounded univariate) Gaussians. If the covariance matrices of each of the Gaussians is the identity matrix, we show that $\widetilde{O}(kd\alpha^2 + kd \log(1/\delta) / \alpha \epsilon)$ samples are sufficient. To prove our results, we design a new technique for privately learning mixture distributions. A class of distributions \mathcal{F} is said to be list-decodable if there is an algorithm that, given "heavily corrupted" samples from $f \in \mathcal{F}$, outputs a list of distributions one of which approximates f . We show that if \mathcal{F} is privately list-decodable then we can learn mixtures of distributions in \mathcal{F} . Finally, we show axis-aligned Gaussian distributions are privately list-decodable, thereby proving mixtures of such distributions are privately learnable.

[Deep Self-Dissimilarities as Powerful Visual Fingerprints](#)

- Idan Kligvasser, Tamar Shaham, Yuval Bahat, Tomer Michaeli
- abstract@[open-review](#): Features extracted from deep layers of classification networks are widely used as image descriptors. Here, we exploit an unexplored property of these features: their internal dissimilarity. While small image patches are known to have similar statistics across image scales, it turns out that the internal distribution of deep features varies distinctively between scales. We show how this deep self dissimilarity (DSD) property can be used as a powerful visual fingerprint. Particularly, we illustrate that full-reference and no-reference image quality measures derived from DSD are highly correlated with human preference. In addition, incorporating DSD as a loss function in training of image restoration networks, leads to results that are at least as photo-realistic as those obtained by GAN based methods, while not requiring adversarial training.

[Invariant Causal Imitation Learning for Generalizable Policies](#)

- Ioana Bica, Daniel Jarrett, Mihaela van der Schaar
- abstract@[open-review](#): Consider learning an imitation policy on the basis of demonstrated behavior from multiple environments, with an eye towards deployment in an unseen environment. Since the observable features from each setting may be different, directly learning individual policies as mappings from features to actions is prone to spurious correlations---and may not generalize well. However, the expert's policy is often a function of a shared latent structure underlying those observable features that is invariant across settings. By leveraging data from multiple environments, we propose Invariant Causal Imitation Learning (ICIL), a novel technique in which we learn a feature representation that is invariant across domains, on the basis of which we learn an imitation policy that matches expert behavior. To cope with transition dynamics mismatch, ICIL learns a shared representation of causal features (for all training environments), that is disentangled from the specific representations of noise variables (for each of those environments). Moreover, to ensure that the learned policy matches the observation distribution of the expert's policy, ICIL estimates the energy of the expert's observations and uses a regularization term that minimizes the imitator policy's next state energy. Experimentally, we compare our methods against several benchmarks in control and healthcare tasks and show its effectiveness in learning imitation policies capable of generalizing to unseen environments.

[CoAtNet: Marrying Convolution and Attention for All Data Sizes](#)

- Zihang Dai, Hanxiao Liu, Quoc V Le, Mingxing Tan
- abstract@[open-review](#): Transformers have attracted increasing interests in computer vision, but they still fall behind state-of-the-art convolutional networks. In this work, we show that while Transformers tend to have larger model capacity, their generalization can be worse than convolutional networks due to the lack of the right inductive bias. To effectively combine the strengths from both architectures, we present CoAtNets(pronounced "coat" nets), a family of hybrid models built from two key insights: (1) depthwise Convolution and self-Attention can be naturally unified via simple relative attention; (2) vertically stacking convolution layers and attention layers in a principled way is surprisingly effective in improving generalization, capacity and efficiency. Experiments show that our CoAtNets achieve state-of-the-art performance under different resource constraints across various datasets: Without extra data, CoAtNet achieves 86.0% ImageNet top-1 accuracy; When pre-trained with 13M images from ImageNet-21K, our CoAtNet achieves 88.56% top-1 accuracy, matching ViT-huge pre-trained with 300M images from JFT-300M while using 23x less data; Notably, when we further scale up CoAtNet with JFT-3B, it achieves 90.88% top-1 accuracy on ImageNet, establishing a new state-of-the-art result.

[Mixed Supervised Object Detection by Transferring Mask Prior and Semantic Similarity](#)

- Yan Liu, Zhijie Zhang, Li Niu, Junjie Chen, Liqing Zhang
- abstract@[open-review](#): Object detection has achieved promising success, but requires large-scale fully-annotated data, which is time-consuming and labor-extensive. Therefore, we consider object detection with mixed supervision, which learns novel object categories using weak annotations with the help of full annotations of existing base object categories. Previous works using mixed supervision mainly learn the class-agnostic objectness from fully-annotated categories, which can be transferred to upgrade the weak annotations to pseudo full annotations for novel categories. In this paper, we further transfer mask prior and semantic similarity to bridge the gap between novel categories and base categories. Specifically, the ability of using mask prior to help detect objects is learned from base categories and transferred to novel categories. Moreover, the semantic similarity between objects learned from base categories is transferred to denoise the pseudo full annotations for novel categories. Experimental results on three benchmark datasets demonstrate the effectiveness of our method over existing methods. Codes are available at <https://github.com/bcmi/TraMaS-Weak-Shot-Object-Detection>.

[Celebrating Diversity in Shared Multi-Agent Reinforcement Learning](#)

- Chenghao Li, Tonghan Wang, Chengjie Wu, Qianchuan Zhao, Jun Yang, Chongjie Zhang
- abstract@[open-review](#): Recently, deep multi-agent reinforcement learning (MARL) has shown the promise to solve complex cooperative tasks. Its success is partly because of parameter sharing among agents. However, such sharing may lead agents to behave similarly and limit their coordination capacity. In this paper, we aim to introduce diversity in both optimization and representation of shared multi-agent reinforcement learning. Specifically, we propose an

information-theoretical regularization to maximize the mutual information between agents' identities and their trajectories, encouraging extensive exploration and diverse individualized behaviors. In representation, we incorporate agent-specific modules in the shared neural network architecture, which are regularized by L1-norm to promote learning sharing among agents while keeping necessary diversity. Empirical results show that our method achieves state-of-the-art performance on Google Research Football and super hard StarCraft II micromanagement tasks.

[Rebounding Bandits for Modeling Satiation Effects](#)

- Liu Leqi, Fatma Kilinc Karzan, Zachary Lipton, Alan Montgomery
- abstract@[open-review](#): Psychological research shows that enjoyment of many goods is subject to satiation, with short-term satisfaction declining after repeated exposures to the same item. Nevertheless, proposed algorithms for powering recommender systems seldom model these dynamics, instead proceeding as though user preferences were fixed in time. In this work, we introduce rebounding bandits, a multi-armed bandit setup, where satiation dynamics are modeled as time-invariant linear dynamical systems. Expected rewards for each arm decline monotonically with consecutive exposures and rebound towards the initial reward whenever that arm is not pulled. Unlike classical bandit algorithms, methods for tackling rebounding bandits must plan ahead and model-based methods rely on estimating the parameters of the satiation dynamics. We characterize the planning problem, showing that the greedy policy is optimal when the arms exhibit identical deterministic dynamics. To address stochastic satiation dynamics with unknown parameters, we propose Explore-Estimate-Plan, an algorithm that pulls arms methodically, estimates the system dynamics, and then plans accordingly.

[Sample Complexity of Tree Search Configuration: Cutting Planes and Beyond](#)

- Maria-Florina F. Balcan, Siddharth Prasad, Tuomas Sandholm, Ellen Vitercik
- abstract@[open-review](#): Cutting-plane methods have enabled remarkable successes in integer programming over the last few decades. State-of-the-art solvers integrate a myriad of cutting-plane techniques to speed up the underlying tree-search algorithm used to find optimal solutions. In this paper we provide sample complexity bounds for cut-selection in branch-and-cut (B&C). Given a training set of integer programs sampled from an application-specific input distribution and a family of cut selection policies, these guarantees bound the number of samples sufficient to ensure that using any policy in the family, the size of the tree B&C builds on average over the training set is close to the expected size of the tree B&C builds. We first bound the sample complexity of learning cutting planes from the canonical family of Chvátal-Gomory cuts. Our bounds handle any number of waves of any number of cuts and are fine tuned to the magnitudes of the constraint coefficients. Next, we prove sample complexity bounds for more sophisticated cut selection policies that use a combination of scoring rules to choose from a family of cuts. Finally, beyond the realm of cutting planes for integer programming, we develop a general abstraction of tree search that captures key components such as node selection and variable selection. For this abstraction, we bound the sample complexity of learning a good policy for building the search tree.

[IQ-Learn: Inverse soft-Q Learning for Imitation](#)

- Divyansh Garg, Shuvam Chakraborty, Chris Cundy, Jiaming Song, Stefano Ermon
- abstract@[open-review](#): In many sequential decision-making problems (e.g., robotics control, game playing, sequential prediction), human or expert data is available containing useful information about the task. However, imitation learning (IL) from a small amount of expert data can be challenging in high-dimensional environments with complex dynamics. Behavioral cloning is a simple method that is widely used due to its simplicity of implementation and stable convergence but doesn't utilize any information involving the environment's dynamics. Many existing methods that exploit dynamics information are difficult to train in practice due to an adversarial optimization process over reward and policy approximators or biased, high variance gradient estimators. We introduce a method for dynamics-aware IL which avoids adversarial training by learning a single Q-function, implicitly representing both reward and policy. On standard benchmarks, the implicitly learned rewards show a high positive correlation with the ground-truth rewards, illustrating our method can also be used for inverse reinforcement learning (IRL). Our method, Inverse soft-Q learning (IQ-Learn) obtains state-of-the-art results in offline and online imitation learning settings, significantly outperforming existing methods both in the number of required environment interactions and scalability in high-dimensional spaces, often by more than 3x.

[Task-Agnostic Undesirable Feature Deactivation Using Out-of-Distribution Data](#)

- Dongmin Park, Hwanjun Song, Minseok Kim, Jae-Gil Lee
- abstract@[open-review](#): A deep neural network (DNN) has achieved great success in many machine learning tasks by virtue of its high expressive power. However, its prediction can be easily biased to undesirable features, which are not essential for solving the target task and are even imperceptible to a human, thereby resulting in poor generalization. Leveraging plenty of undesirable features in out-of-distribution (OOD) examples has emerged as a potential solution for de-biasing such features, and a recent study shows that softmax-level calibration of OOD examples can successfully remove the contribution of undesirable features to the last fully-connected layer of a classifier. However, its applicability is confined to the classification task, and its impact on a DNN feature extractor is not properly investigated. In this paper, we propose Taufe, a novel regularizer that deactivates many undesirable features using OOD examples in the feature extraction layer and thus removes the dependency on the task-specific softmax layer. To show the task-agnostic nature of Taufe, we rigorously validate its performance on three tasks, classification, regression, and a mix of them, on CIFAR-10, CIFAR-100, ImageNet, CUB200, and CAR datasets. The results demonstrate that Taufe consistently outperforms the state-of-the-art method as well as the baselines without regularization.

[Private Non-smooth ERM and SCO in Subquadratic Steps](#)

- Janardhan Kulkarni, Yin Tat Lee, Daogao Liu
- abstract@[open-review](#): We study the differentially private Empirical Risk Minimization (ERM) and Stochastic Convex Optimization (SCO) problems for non-smooth convex functions. We get a (nearly) optimal bound on the excess empirical risk for ERM with $\$O(\frac{N^{3/2}}{d^{1/8}} + \frac{N^2}{d})$ gradient queries, which is achieved with the help of subsampling and smoothing the function via convolution. Combining this result with the iterative localization technique of Feldman et al. [\cite{fkt20}](#), we achieve the optimal excess population loss for the SCO problem with $\$O(\min\{N^{5/4}d^{1/8}, \frac{N^{3/2}}{d^{1/8}}\})$ gradient queries. Our work makes progress towards resolving a question raised by Bassily et al. [\cite{bfgt20}](#), giving first algorithms for private SCO with subquadratic steps. In a concurrent work, Asi et al. [\cite{afkt21}](#) gave other algorithms for private ERM and SCO with subquadratic steps.

[Towards Instance-Optimal Offline Reinforcement Learning with Pessimism](#)

- Ming Yin, Yu-Xiang Wang
- abstract@[open-review](#): We study the offline reinforcement learning (offline RL) problem, where the goal is to learn a reward-maximizing policy in an unknown Markov Decision Process (MDP) using the data coming from a policy μ . In particular, we consider the sample complexity problems of offline RL for the finite horizon MDPs. Prior works derive the information-theoretical lower bounds based on different data-coverage assumptions and their upper bounds are expressed by the covering coefficients which lack the explicit characterization of system quantities. In this work, we analyze the Adaptive Pessimistic Value Iteration (APVI) algorithm and derive the suboptimality upper bound that nearly matches $\left[\sum_{h=1}^H \sum_{s,h,a,h} d^{\pi^*} h(s,h,a,h) \sqrt{\frac{\text{Var}\{P_{s,h,a,h}\}(V^{*\star}_{h+1} + r_h)}{d^{\mu_h(s,h,a,h)}}}\right] \sqrt{\frac{1}{n}}$. We also prove an information-theoretical lower bound to show this quantity is required under the weak assumption that $d^{\mu_h(s,h,a,h)} > 0$ if $d^{\pi^*} h(s,h,a,h) > 0$. Here π^* is a optimal policy, μ is the behavior policy and $d(s,h,a)$ is the marginal state-action probability. We call this adaptive bound the intrinsic offline reinforcement learning bound since it

directly implies all the existing optimal results: minimax rate under uniform data-coverage assumption, horizon-free setting, single policy concentrability, and the tight problem-dependent results. Later, we extend the result to the \emph{assumption-free} regime (where we make no assumption on μ) and obtain the assumption-free intrinsic bound. Due to its generic form, we believe the intrinsic bound could help illuminate what makes a specific problem hard and reveal the fundamental challenges in offline RL.

[Speedy Performance Estimation for Neural Architecture Search](#)

- Robin Ru, Clare Lyle, Lisa Schut, Miroslav Fil, Mark van der Wilk, Yarin Gal
- abstract@[open-review](#): Reliable yet efficient evaluation of generalisation performance of a proposed architecture is crucial to the success of neural architecture search (NAS). Traditional approaches face a variety of limitations: training each architecture to completion is prohibitively expensive, early stopped validation accuracy may correlate poorly with fully trained performance, and model-based estimators require large training sets. We instead propose to estimate the final test performance based on a simple measure of training speed. Our estimator is theoretically motivated by the connection between generalisation and training speed, and is also inspired by the reformulation of a PAC-Bayes bound under the Bayesian setting. Our model-free estimator is simple, efficient, and cheap to implement, and does not require hyperparameter-tuning or surrogate training before deployment. We demonstrate on various NAS search spaces that our estimator consistently outperforms other alternatives in achieving better correlation with the true test performance rankings. We further show that our estimator can be easily incorporated into both query-based and one-shot NAS methods to improve the speed or quality of the search.

[How Tight Can PAC-Bayes be in the Small Data Regime?](#)

- Andrew Foong, Wessel Bruinsma, David Burt, Richard Turner
- abstract@[open-review](#): In this paper, we investigate the question: *Given a small number of datapoints, for example $N = 30$, how tight can PAC-Bayes and test set bounds be made?* For such small datasets, test set bounds adversely affect generalisation performance by withholding data from the training procedure. In this setting, PAC-Bayes bounds are especially attractive, due to their ability to use all the data to simultaneously learn a posterior and bound its generalisation risk. We focus on the case of i.i.d. data with a bounded loss and consider the generic PAC-Bayes theorem of Germain et al. While their theorem is known to recover many existing PAC-Bayes bounds, it is unclear what the tightest bound derivable from their framework is. For a fixed learning algorithm and dataset, we show that the tightest possible bound coincides with a bound considered by Catoni; and, in the more natural case of distributions over datasets, we establish a lower bound on the best bound achievable in expectation. Interestingly, this lower bound recovers the Chernoff test set bound if the posterior is equal to the prior. Moreover, to illustrate how tight these bounds can be, we study synthetic one-dimensional classification tasks in which it is feasible to meta-learn both the prior and the form of the bound to numerically optimise for the tightest bounds possible. We find that in this simple, controlled scenario, PAC-Bayes bounds are competitive with comparable, commonly used Chernoff test set bounds. However, the sharpest test set bounds still lead to better guarantees on the generalisation error than the PAC-Bayes bounds we consider.

[Deep Synoptic Monte-Carlo Planning in Reconnaissance Blind Chess](#)

- Gregory Clark
- abstract@[open-review](#): This paper introduces deep synoptic Monte Carlo planning (DSMCP) for large imperfect information games. The algorithm constructs a belief state with an unweighted particle filter and plans via playouts that start at samples drawn from the belief state. The algorithm accounts for uncertainty by performing inference on "synopses," a novel stochastic abstraction of information states. DSMCP is the basis of the program Penumbra, which won the official 2020 reconnaissance blind chess competition versus 33 other programs. This paper also evaluates algorithm variants that incorporate caution, paranoia, and a novel bandit algorithm. Furthermore, it audits the synopsis features used in Penumbra with per-bit saliency statistics.

[Dynamic Analysis of Higher-Order Coordination in Neuronal Assemblies via De-Sparsified Orthogonal Matching Pursuit](#)

- Shoutik Mukherjee, Behtash Babadi
- abstract@[open-review](#): Coordinated ensemble spiking activity is widely observable in neural recordings and central in the study of population codes, with hypothesized roles including robust stimulus representation, interareal communication of neural information, and learning and memory formation. Model-free measures of synchrony characterize the coherence of pairwise activity, but not higher-order interactions; this limitation is transcended by statistical models of ensemble spiking activity. However, existing model-based analyses often impose assumptions about the relevance of higher-order interactions and require multiple repeated trials in order to characterize dynamics in the correlational structure of ensemble activity. To address these shortcomings, we propose an adaptive greedy filtering algorithm based on a discretized mark point-process model of ensemble spiking and a corresponding precise statistical inference framework to identify significant coordinated higher-order spiking activity. In the course of developing the statistical inference procedures, we also show that confidence intervals can be constructed for greedily estimated parameters. We demonstrate the utility of our proposed methods on simulated neuronal assemblies. Applied to multi-electrode recordings of human cortical ensembles, our proposed methods provide new insights into the dynamics underlying localized population activity during transitions between brain states.

[Efficient Training of Retrieval Models using Negative Cache](#)

- Erik Lindgren, Sashank Reddi, Ruiqi Guo, Sanjiv Kumar
- abstract@[open-review](#): Factorized models, such as two tower neural network models, are widely used for scoring (query, document) pairs in information retrieval tasks. These models are typically trained by optimizing the model parameters to score relevant positive" pairs higher than the irrelevantnegative" ones. While a large set of negatives typically improves the model performance, limited computation and memory budgets place constraints on the number of negatives used during training. In this paper, we develop a novel negative sampling technique for accelerating training with softmax cross-entropy loss. By using cached (possibly stale) item embeddings, our technique enables training with a large pool of negatives with reduced memory and computation. We also develop a streaming variant of our algorithm geared towards very large datasets. Furthermore, we establish a theoretical basis for our approach by showing that updating a very small fraction of the cache at each iteration can still ensure fast convergence. Finally, we experimentally validate our approach and show that it is efficient and compares favorably with more complex, state-of-the-art approaches.

[Understanding Partial Multi-Label Learning via Mutual Information](#)

- Xiuwen Gong, Dong Yuan, Wei Bao
- abstract@[open-review](#): To deal with ambiguities in partial multilabel learning (PML), state-of-the-art methods perform disambiguation by identifying ground-truth labels directly. However, there is an essential question: "Can the ground-truth labels be identified precisely?". If yes, "How can the ground-truth labels be found?". This paper provides affirmative answers to these questions. Instead of adopting hand-made heuristic strategy, we propose a novel Mutual Information Label Identification for Partial Multilabel Learning (MILI-PML), which is derived from a clear probabilistic formulation and could be easily interpreted theoretically from the mutual information perspective, as well as naturally incorporates the feature/label relevancy considerations. Extensive experiments on synthetic and real-world datasets clearly demonstrate the superiorities of the proposed MILI-PML.

[Environment Generation for Zero-Shot Compositional Reinforcement Learning](#)

- Izzeddin Gur, Natasha Jaques, Yingjie Miao, Jongwook Choi, Manoj Tiwari, Honglak Lee, Aleksandra Faust

- abstract@[open-review](#): Many real-world problems are compositional – solving them requires completing interdependent sub-tasks, either in series or in parallel, that can be represented as a dependency graph. Deep reinforcement learning (RL) agents often struggle to learn such complex tasks due to the long time horizons and sparse rewards. To address this problem, we present Compositional Design of Environments (CoDE), which trains a Generator agent to automatically build a series of compositional tasks tailored to the RL agent’s current skill level. This automatic curriculum not only enables the agent to learn more complex tasks than it could have otherwise, but also selects tasks where the agent’s performance is weak, enhancing its robustness and ability to generalize zero-shot to unseen tasks at test-time. We analyze why current environment generation techniques are insufficient for the problem of generating compositional tasks, and propose a new algorithm that addresses these issues. Our results assess learning and generalization across multiple compositional tasks, including the real-world problem of learning to navigate and interact with web pages. We learn to generate environments composed of multiple pages or rooms, and train RL agents capable of completing wide-range of complex tasks in those environments. We contribute two new benchmark frameworks for generating compositional tasks, compositional MiniGrid and gMiniWoB for web navigation. CoDE yields 4x higher success rate than the strongest baseline, and demonstrates strong performance of real websites learned on 3500 primitive tasks.

[Optimizing Conditional Value-At-Risk of Black-Box Functions](#)

- Quoc Phong Nguyen, Zhongxiang Dai, Bryan Kian Hsiang Low, Patrick Jaillet
- abstract@[open-review](#): This paper presents two Bayesian optimization (BO) algorithms with theoretical performance guarantee to maximize the conditional value-at-risk (CVaR) of a black-box function: CV-UCB and CV-TS which are based on the well-established principle of optimism in the face of uncertainty and Thompson sampling, respectively. To achieve this, we develop an upper confidence bound of CVaR and prove the no-regret guarantee of CV-UCB by utilizing an interesting connection between CVaR and value-at-risk (VaR). For CV-TS, though it is straightforwardly performed with Thompson sampling, bounding its Bayesian regret is non-trivial because it requires a tail expectation bound for the distribution of CVaR of a black-box function, which has not been shown in the literature. The performances of both CV-UCB and CV-TS are empirically evaluated in optimizing CVaR of synthetic benchmark functions and simulated real-world optimization problems.

[E\(n\) Equivariant Normalizing Flows](#)

- Victor Garcia Satorras, Emiel Hoogeboom, Fabian Fuchs, Ingmar Posner, Max Welling
- abstract@[open-review](#): This paper introduces a generative model equivariant to Euclidean symmetries: E(n) Equivariant Normalizing Flows (E-NFs). To construct E-NFs, we take the discriminative E(n) graph neural networks and integrate them as a differential equation to obtain an invertible equivariant function: a continuous-time normalizing flow. We demonstrate that E-NFs considerably outperform baselines and existing methods from the literature on particle systems such as DW4 and LJ13, and on molecules from QM9 in terms of log-likelihood. To the best of our knowledge, this is the first flow that jointly generates molecule features and positions in 3D.

[Revitalizing CNN Attention via Transformers in Self-Supervised Visual Representation Learning](#)

- Chongjian GE, Youwei Liang, YIBING SONG, Jianbo Jiao, Jue Wang, Ping Luo
- abstract@[open-review](#): Studies on self-supervised visual representation learning (SSL) improve encoder backbones to discriminate training samples without labels. While CNN encoders via SSL achieve comparable recognition performance to those via supervised learning, their network attention is under-explored for further improvement. Motivated by the transformers that explore visual attention effectively in recognition scenarios, we propose a CNN Attention REvitalization (CARE) framework to train attentive CNN encoders guided by transformers in SSL. The proposed CARE framework consists of a CNN stream (C-stream) and a transformer stream (T-stream), where each stream contains two branches. C-stream follows an existing SSL framework with two CNN encoders, two projectors, and a predictor. T-stream contains two transformers, two projectors, and a predictor. T-stream connects to CNN encoders and is in parallel to the remaining C-Stream. During training, we perform SSL in both streams simultaneously and use the T-stream output to supervise C-stream. The features from CNN encoders are modulated in T-stream for visual attention enhancement and become suitable for the SSL scenario. We use these modulated features to supervise C-stream for learning attentive CNN encoders. To this end, we revitalize CNN attention by using transformers as guidance. Experiments on several standard visual recognition benchmarks, including image classification, object detection, and semantic segmentation, show that the proposed CARE framework improves CNN encoder backbones to the state-of-the-art performance.

[A Critical Look at the Consistency of Causal Estimation with Deep Latent Variable Models](#)

- Severi Rissanen, Pekka Marttinen
- abstract@[open-review](#): Using deep latent variable models in causal inference has attracted considerable interest recently, but an essential open question is their ability to yield consistent causal estimates. While they have demonstrated promising results and theory exists on some simple model formulations, we also know that causal effects are not even identifiable in general with latent variables. We investigate this gap between theory and empirical results with analytical considerations and extensive experiments under multiple synthetic and real-world data sets, using the causal effect variational autoencoder (CEVAE) as a case study. While CEVAE seems to work reliably under some simple scenarios, it does not estimate the causal effect correctly with a misspecified latent variable or a complex data distribution, as opposed to its original motivation. Hence, our results show that more attention should be paid to ensuring the correctness of causal estimates with deep latent variable models.

[Improving Robustness using Generated Data](#)

- Sven Gowal, Sylvestre-Alvise Rebuffi, Olivia Wiles, Florian Stimberg, Dan Andrei Calian, Timothy A Mann
- abstract@[open-review](#): Recent work argues that robust training requires substantially larger datasets than those required for standard classification. On CIFAR-10 and CIFAR-100, this translates into a sizable robust-accuracy gap between models trained solely on data from the original training set and those trained with additional data extracted from the "80 Million Tiny Images" dataset (TI-80M). In this paper, we explore how generative models trained solely on the original training set can be leveraged to artificially increase the size of the original training set and improve adversarial robustness to $\|\cdot\|_p$ norm-bounded perturbations. We identify the sufficient conditions under which incorporating additional generated data can improve robustness, and demonstrate that it is possible to significantly reduce the robust-accuracy gap to models trained with additional real data. Surprisingly, we even show that even the addition of non-realistic random data (generated by Gaussian sampling) can improve robustness. We evaluate our approach on CIFAR-10, CIFAR-100, SVHN and TinyImageNet against $\|\cdot\|_\infty$ and $\|\cdot\|_2$ norm-bounded perturbations of size $\|\epsilon\| = 8/255$ and $\|\epsilon\| = 128/255$, respectively. We show large absolute improvements in robust accuracy compared to previous state-of-the-art methods. Against $\|\cdot\|_\infty$ norm-bounded perturbations of size $\|\epsilon\| = 8/255$, our models achieve 66.10% and 33.49% robust accuracy on CIFAR-10 and CIFAR-100, respectively (improving upon the state-of-the-art by +8.96% and +3.29%). Against $\|\cdot\|_2$ norm-bounded perturbations of size $\|\epsilon\| = 128/255$, our model achieves 78.31% on CIFAR-10 (+3.81%). These results beat most prior works that use external data.

[An Analysis of Constant Step Size SGD in the Non-convex Regime: Asymptotic Normality and Bias](#)

- Lu Yu, Krishnakumar Balasubramanian, Stanislav Volgshev, Murat A. Erdogdu
- abstract@[open-review](#): Structured non-convex learning problems, for which critical points have favorable statistical properties, arise frequently in statistical machine learning. Algorithmic convergence and statistical estimation rates are well-understood for such problems. However, quantifying the uncertainty associated with the underlying training algorithm is not well-studied in the non-convex setting. In order to address this shortcoming, in this work, we establish an asymptotic normality result for the constant step size stochastic gradient descent (SGD) algorithm---a widely used algorithm in practice. Specifically, based on the relationship between SGD and Markov Chains [DDB19], we show that the average of SGD iterates is asymptotically

normally distributed around the expected value of their unique invariant distribution, as long as the non-convex and non-smooth objective function satisfies a dissipativity property. We also characterize the bias between this expected value and the critical points of the objective function under various local regularity conditions. Together, the above two results could be leveraged to construct confidence intervals for non-convex problems that are trained using the SGD algorithm.

[Learning to Learn Graph Topologies](#)

- Xingyue Pu, Tianyue Cao, Xiaoyun Zhang, Xiaowen Dong, Siheng Chen
- abstract@[open-review](#): Learning a graph topology to reveal the underlying relationship between data entities plays an important role in various machine learning and data analysis tasks. Under the assumption that structured data vary smoothly over a graph, the problem can be formulated as a regularised convex optimisation over a positive semidefinite cone and solved by iterative algorithms. Classic methods require an explicit convex function to reflect generic topological priors, e.g. the $\|\cdot\|_1$ penalty for enforcing sparsity, which limits the flexibility and expressiveness in learning rich topological structures. We propose to learn a mapping from node data to the graph structure based on the idea of learning to optimise (L2O). Specifically, our model first unrolls an iterative primal-dual splitting algorithm into a neural network. The key structural proximal projection is replaced with a variational autoencoder that refines the estimated graph with enhanced topological properties. The model is trained in an end-to-end fashion with pairs of node data and graph samples. Experiments on both synthetic and real-world data demonstrate that our model is more efficient than classic iterative algorithms in learning a graph with specific topological properties.

[Invertible Tabular GANs: Killing Two Birds with One Stone for Tabular Data Synthesis](#)

- JAEHOON LEE, Jihyeon Hyeong, Jinsung Jeon, Noseong Park, Jihoon Cho
- abstract@[open-review](#): Tabular data synthesis has received wide attention in the literature. This is because available data is often limited, incomplete, or cannot be obtained easily, and data privacy is becoming increasingly important. In this work, we present a generalized GAN framework for tabular synthesis, which combines the adversarial training of GANs and the negative log-density regularization of invertible neural networks. The proposed framework can be used for two distinctive objectives. First, we can further improve the synthesis quality, by decreasing the negative log-density of real records in the process of adversarial training. On the other hand, by increasing the negative log-density of real records, realistic fake records can be synthesized in a way that they are not too much close to real records and reduce the chance of potential information leakage. We conduct experiments with real-world datasets for classification, regression, and privacy attacks. In general, the proposed method demonstrates the best synthesis quality (in terms of task-oriented evaluation metrics, e.g., F1) when decreasing the negative log-density during the adversarial training. If increasing the negative log-density, our experimental results show that the distance between real and fake records increases, enhancing robustness against privacy attacks.

[Reducing Collision Checking for Sampling-Based Motion Planning Using Graph Neural Networks](#)

- Chennng Yu, Sicun Gao
- abstract@[open-review](#): Sampling-based motion planning is a popular approach in robotics for finding paths in continuous configuration spaces. Checking collision with obstacles is the major computational bottleneck in this process. We propose new learning-based methods for reducing collision checking to accelerate motion planning by training graph neural networks (GNNs) that perform path exploration and path smoothing. Given random geometric graphs (RGGs) generated from batch sampling, the path exploration component iteratively predicts collision-free edges to prioritize their exploration. The path smoothing component then optimizes paths obtained from the exploration stage. The methods benefit from the ability of GNNs of capturing geometric patterns from RGGs through batch sampling and generalize better to unseen environments. Experimental results show that the learned components can significantly reduce collision checking and improve overall planning efficiency in challenging high-dimensional motion planning tasks.

[Sample Complexity Bounds for Active Ranking from Multi-wise Comparisons](#)

- Wenbo Ren, Jia Liu, Ness Shroff
- abstract@[open-review](#): We study the sample complexity (i.e., the number of comparisons needed) bounds for actively ranking a set of n items from multi-wise comparisons. Here, a multi-wise comparison takes m items as input and returns a (noisy) result about the best item (the winner feedback) or the order of these items (the full-ranking feedback). We consider two basic ranking problems: top- k items selection and full ranking. Unlike previous works that study ranking from multi-wise comparisons, in this paper, we do not require any parametric model or assumption and work on the fundamental setting where each comparison returns the correct result with probability 1 or a certain probability larger than $\frac{1}{2}$. This paper helps understand whether and to what degree utilizing multi-wise comparisons can reduce the sample complexity for the ranking problems compared to ranking from pairwise comparisons. Specifically, under the winner feedback setting, one can reduce the sample complexity for top- k selection up to an m factor and that for full ranking up to a $\log m$ factor. Under the full-ranking feedback setting, one can reduce the sample complexity for top- k selection up to an m factor and that for full ranking up to an $m \log m$ factor. We also conduct numerical simulations to confirm our theoretical results.

[Efficient Bayesian network structure learning via local Markov boundary search](#)

- Ming Gao, Bryon Aragam
- abstract@[open-review](#): We analyze the complexity of learning directed acyclic graphical models from observational data in general settings without specific distributional assumptions. Our approach is information-theoretic and uses a local Markov boundary search procedure in order to recursively construct ancestral sets in the underlying graphical model. Perhaps surprisingly, we show that for certain graph ensembles, a simple forward greedy search algorithm (i.e. without a backward pruning phase) suffices to learn the Markov boundary of each node. This substantially improves the sample complexity, which we show is at most polynomial in the number of nodes. This is then applied to learn the entire graph under a novel identifiability condition that generalizes existing conditions from the literature. As a matter of independent interest, we establish finite-sample guarantees for the problem of recovering Markov boundaries from data. Moreover, we apply our results to the special case of polytrees, for which the assumptions simplify, and provide explicit conditions under which polytrees are identifiable and learnable in polynomial time. We further illustrate the performance of the algorithm, which is easy to implement, in a simulation study. Our approach is general, works for discrete or continuous distributions without distributional assumptions, and as such sheds light on the minimal assumptions required to efficiently learn the structure of directed graphical models from data.

[Learning Dynamic Graph Representation of Brain Connectome with Spatio-Temporal Attention](#)

- Byung-Hoon Kim, Jong Chul Ye, Jae-Jin Kim
- abstract@[open-review](#): Functional connectivity (FC) between regions of the brain can be assessed by the degree of temporal correlation measured with functional neuroimaging modalities. Based on the fact that these connectivities build a network, graph-based approaches for analyzing the brain connectome have provided insights into the functions of the human brain. The development of graph neural networks (GNNs) capable of learning representation from graph structured data has led to increased interest in learning the graph representation of the brain connectome. Although recent attempts to apply GNN to the FC network have shown promising results, there is still a common limitation that they usually do not incorporate the dynamic characteristics of the FC network which fluctuates over time. In addition, a few studies that have attempted to use dynamic FC as an input for the GNN reported a reduction in performance compared to static FC methods, and did not provide temporal explainability. Here, we propose STAGIN, a method for learning dynamic graph representation of the brain connectome with spatio-temporal attention. Specifically, a temporal sequence of brain graphs is input to the STAGIN to obtain the dynamic graph representation, while novel READOUT functions and the Transformer encoder provide spatial

and temporal explainability with attention, respectively. Experiments on the HCP-Rest and the HCP-Task datasets demonstrate exceptional performance of our proposed method. Analysis of the spatio-temporal attention also provide concurrent interpretation with the neuroscientific knowledge, which further validates our method. Code is available at <https://github.com/egyptdj/stagin>

Understanding the Generalization Benefit of Model Invariance from a Data Perspective

- Sicheng Zhu, Bang An, Furong Huang
- abstract@[open-review](#): Machine learning models that are developed to be invariant under certain types of data transformations have shown improved generalization in practice. However, a principled understanding of why invariance benefits generalization is limited. Given a dataset, there is often no principled way to select "suitable" data transformations under which model invariance guarantees better generalization. This paper studies the generalization benefit of model invariance by introducing the sample cover induced by transformations, i.e., a representative subset of a dataset that can approximately recover the whole dataset using transformations. For any data transformations, we provide refined generalization bounds for invariant models based on the sample cover. We also characterize the "suitability" of a set of data transformations by the sample covering number induced by transformations, i.e., the smallest size of its induced sample covers. We show that we may tighten the generalization bounds for "suitable" transformations that have a small sample covering number. In addition, our proposed sample covering number can be empirically evaluated and thus provides a guidance for selecting transformations to develop model invariance for better generalization. In experiments on multiple datasets, we evaluate sample covering numbers for some commonly used transformations and show that the smaller sample covering number for a set of transformations (e.g., the 3D-view transformation) indicates a smaller gap between the test and training error for invariant models, which verifies our propositions.

Improved Variance-Aware Confidence Sets for Linear Bandits and Linear Mixture MDP

- Zihan Zhang, Jiaqi Yang, Xiangyang Ji, Simon S. Du
- abstract@[open-review](#): This paper presents new \emph{variance-aware} confidence sets for linear bandits and linear mixture Markov Decision Processes (MDPs). With the new confidence sets, we obtain the follow regret bounds: For linear bandits, we obtain an $\widetilde{O}(\mathrm{poly}(d)\sqrt{1 + \sum_{k=1}^K \sigma_k^2})$ data-dependent regret bound, where d is the feature dimension, K is the number of rounds, and σ_k^2 is the \emph{unknown} variance of the reward at the k -th round. This is the first regret bound that only scales with the variance and the dimension but \emph{no explicit polynomial dependency on K }. When variances are small, this bound can be significantly smaller than the $\widetilde{O}(d\sqrt{K})$ worst-case regret bound. For linear mixture MDPs, we obtain an $\widetilde{O}(\mathrm{poly}(d, \log H)\sqrt{K})$ regret bound, where d is the number of base models, K is the number of episodes, and H is the planning horizon. This is the first regret bound that only scales \emph{logarithmically} with H in the reinforcement learning with linear function approximation setting, thus \emph{exponentially improving} existing results, and resolving an open problem in \citet{zhou2020nearly}. We develop three technical ideas that may be of independent interest: 1) applications of the peeling technique to both the input norm and the variance magnitude, 2) a recursion-based estimator for the variance, and 3) a new convex potential lemma that generalizes the seminal elliptical potential lemma.

How Should Pre-Trained Language Models Be Fine-Tuned Towards Adversarial Robustness?

- Xinshuai Dong, Anh Tuan Luu, Min Lin, Shuicheng Yan, Hanwang Zhang
- abstract@[open-review](#): The fine-tuning of pre-trained language models has a great success in many NLP fields. Yet, it is strikingly vulnerable to adversarial examples, e.g., word substitution attacks using only synonyms can easily fool a BERT-based sentiment analysis model. In this paper, we demonstrate that adversarial training, the prevalent defense technique, does not directly fit a conventional fine-tuning scenario, because it suffers severely from catastrophic forgetting: failing to retain the generic and robust linguistic features that have already been captured by the pre-trained model. In this light, we propose Robust Informative Fine-Tuning (RIFT), a novel adversarial fine-tuning method from an information-theoretical perspective. In particular, RIFT encourages an objective model to retain the features learned from the pre-trained model throughout the entire fine-tuning process, whereas a conventional one only uses the pre-trained weights for initialization. Experimental results show that RIFT consistently outperforms the state-of-the-arts on two popular NLP tasks: sentiment analysis and natural language inference, under different attacks across various pre-trained language models.

Recursive Bayesian Networks: Generalising and Unifying Probabilistic Context-Free Grammars and Dynamic Bayesian Networks

- Robert Lieck, Martin Rohrmeier
- abstract@[open-review](#): Probabilistic context-free grammars (PCFGs) and dynamic Bayesian networks (DBNs) are widely used sequence models with complementary strengths and limitations. While PCFGs allow for nested hierarchical dependencies (tree structures), their latent variables (non-terminal symbols) have to be discrete. In contrast, DBNs allow for continuous latent variables, but the dependencies are strictly sequential (chain structure). Therefore, neither can be applied if the latent variables are assumed to be continuous and also to have a nested hierarchical dependency structure. In this paper, we present Recursive Bayesian Networks (RBMs), which generalise and unify PCFGs and DBNs, combining their strengths and containing both as special cases. RBMs define a joint distribution over tree-structured Bayesian networks with discrete or continuous latent variables. The main challenge lies in performing joint inference over the exponential number of possible structures and the continuous variables. We provide two solutions: 1) For arbitrary RBMs, we generalise inside and outside probabilities from PCFGs to the mixed discrete-continuous case, which allows for maximum posterior estimates of the continuous latent variables via gradient descent, while marginalising over network structures. 2) For Gaussian RBMs, we additionally derive an analytic approximation of the marginal data likelihood (evidence) and marginal posterior distribution, allowing for robust parameter optimisation and Bayesian inference. The capacity and diverse applications of RBMs are illustrated on two examples: In a quantitative evaluation on synthetic data, we demonstrate and discuss the advantage of RBMs for segmentation and tree induction from noisy sequences, compared to change point detection and hierarchical clustering. In an application to musical data, we approach the unsolved problem of hierarchical music analysis from the raw note level and compare our results to expert annotations.

EF21: A New, Simpler, Theoretically Better, and Practically Faster Error Feedback

- Peter Richtarik, Igor Sokolov, Ilyas Fatkhullin
- abstract@[open-review](#): Error feedback (EF), also known as error compensation, is an immensely popular convergence stabilization mechanism in the context of distributed training of supervised machine learning models enhanced by the use of contractive communication compression mechanisms, such as Top-\$k\$. First proposed by Seide et al [2014] as a heuristic, EF resisted any theoretical understanding until recently [Stich et al., 2018, Alistarh et al., 2018]. While these early breakthroughs were followed by a steady stream of works offering various improvements and generalizations, the current theoretical understanding of EF is still very limited. Indeed, to the best of our knowledge, all existing analyses either i) apply to the single node setting only, ii) rely on very strong and often unreasonable assumptions, such as global boundedness of the gradients, or iterate-dependent assumptions that cannot be checked a-priori and may not hold in practice, or iii) circumvent these issues via the introduction of additional unbiased compressors, which increase the communication cost. In this work we fix all these deficiencies by proposing and analyzing a new EF mechanism, which we call EF21, which consistently and substantially outperforms EF in practice. Moreover, our theoretical analysis relies on standard assumptions only, works in the distributed heterogeneous data setting, and leads to better and more meaningful rates. In particular, we prove that EF21 enjoys a fast $\mathcal{O}(1/T)$ convergence rate for smooth nonconvex problems, beating the previous bound of $\mathcal{O}(1/T^{2/3})$, which was shown under a strong bounded gradients assumption. We further improve this to a fast linear rate for Polyak-Lojasiewicz functions, which is the first linear convergence result for an error feedback method not relying on unbiased compressors. Since EF has a large number of applications where it reigns supreme, we believe that our 2021 variant, EF21, will have a large impact on the practice of communication efficient distributed learning.

[Mixture weights optimisation for Alpha-Divergence Variational Inference](#)

- Kamélia Daudel, randal douc
- abstract@[open-review](#): This paper focuses on α -divergence minimisation methods for Variational Inference. More precisely, we are interested in algorithms optimising the mixture weights of any given mixture model, without any information on the underlying distribution of its mixture components parameters. The Power Descent, defined for all $\alpha \neq 1$, is one such algorithm and we establish in our work the full proof of its convergence towards the optimal mixture weights when $\alpha < 1$. Since the α -divergence recovers the widely-used forward Kullback-Leibler when $\alpha \rightarrow 1$, we then extend the Power Descent to the case $\alpha = 1$ and show that we obtain an Entropic Mirror Descent. This leads us to investigate the link between Power Descent and Entropic Mirror Descent: first-order approximations allow us to introduce the R ϵ nyi Descent, a novel algorithm for which we prove an $O(1/N)$ convergence rate. Lastly, we compare numerically the behavior of the unbiased Power Descent and of the biased R ϵ nyi Descent and we discuss the potential advantages of one algorithm over the other.

[Instance-dependent Label-noise Learning under a Structural Causal Model](#)

- Yu Yao, Tongliang Liu, Mingming Gong, Bo Han, Gang Niu, Kun Zhang
- abstract@[open-review](#): Label noise generally degenerates the performance of deep learning algorithms because deep neural networks easily overfit label errors. Let X and Y denote the instance and clean label, respectively. When Y is a cause of X , according to which many datasets have been constructed, e.g., SVHN and CIFAR, the distributions of $P(X)$ and $P(Y|X)$ are generally entangled. This means that the unsupervised instances are helpful to learn the classifier and thus reduce the side effect of label noise. However, it remains elusive on how to exploit the causal information to handle the label-noise problem. We propose to model and make use of the causal process in order to correct the label-noise effect. Empirically, the proposed method outperforms all state-of-the-art methods on both synthetic and real-world label-noise datasets.

[Combining Human Predictions with Model Probabilities via Confusion Matrices and Calibration](#)

- Gavin Kerrigan, Padhraic Smyth, Mark Steyvers
- abstract@[open-review](#): An increasingly common use case for machine learning models is augmenting the abilities of human decision makers. For classification tasks where neither the human nor model are perfectly accurate, a key step in obtaining high performance is combining their individual predictions in a manner that leverages their relative strengths. In this work, we develop a set of algorithms that combine the probabilistic output of a model with the class-level output of a human. We show theoretically that the accuracy of our combination model is driven not only by the individual human and model accuracies, but also by the model's confidence. Empirical results on image classification with CIFAR-10 and a subset of ImageNet demonstrate that such human-model combinations consistently have higher accuracies than the model or human alone, and that the parameters of the combination method can be estimated effectively with as few as ten labeled datapoints.

[LeadCache: Regret-Optimal Caching in Networks](#)

- Debjit Paria, Abhishek Sinha
- abstract@[open-review](#): We consider an online prediction problem in the context of network caching. Assume that multiple users are connected to several caches via a bipartite network. At any time slot, each user may request an arbitrary file chosen from a large catalog. A user's request at a slot is met if the requested file is cached in at least one of the caches connected to the user. Our objective is to predict, prefetch, and optimally distribute the files on the caches at each slot to maximize the total number of cache hits. The problem is non-trivial due to the non-convex and non-smooth nature of the objective function. In this paper, we propose LeadCache - an efficient online caching policy based on the Follow-the-Perturbed-Leader paradigm. We show that LeadCache is regret-optimal up to a factor of $\tilde{O}(n^{3/8})$, where n is the number of users. We design two efficient implementations of the LeadCache policy, one based on Pipage rounding and the other based on Madow's sampling, each of which makes precisely one call to an LP-solver per iteration. Furthermore, with a Strong-Law-type assumption, we show that the total number of file fetches under LeadCache remains almost surely finite over an infinite horizon. Finally, we derive an approximately tight regret lower bound using results from graph coloring. We conclude that the learning-based LeadCache policy decisively outperforms the state-of-the-art caching policies both theoretically and empirically.

[Probabilistic Attention for Interactive Segmentation](#)

- Prasad Gabbur, Manjot Bilkhu, Javier Movellan
- abstract@[open-review](#): We provide a probabilistic interpretation of attention and show that the standard dot-product attention in transformers is a special case of Maximum A Posteriori (MAP) inference. The proposed approach suggests the use of Expectation Maximization algorithms for on-line adaptation of key and value model parameters. This approach is useful for cases in which external agents, e.g., annotators, provide inference-time information about the correct values of some tokens, e.g., the semantic category of some pixels, and we need for this new information to propagate to other tokens in a principled manner. We illustrate the approach on an interactive semantic segmentation task in which annotators and models collaborate online to improve annotation efficiency. Using standard benchmarks, we observe that key adaptation boosts model performance ($\sim 10\%$ mIoU) in the low feedback regime and value propagation improves model responsiveness in the high feedback regime. A PyTorch layer implementation of our probabilistic attention model is available here: <https://github.com/apple/ml-probabilistic-attention>.

[Influence Patterns for Explaining Information Flow in BERT](#)

- Kaiji Lu, Zifan Wang, Piotr Mardziel, Anupam Datta
- abstract@[open-review](#): While attention is all you need may be proving true, we do not know why: attention-based transformer models such as BERT are superior but how information flows from input tokens to output predictions are unclear. We introduce influence patterns, abstractions of sets of paths through a transformer model. Patterns quantify and localize the flow of information to paths passing through a sequence of model nodes. Experimentally, we find that significant portion of information flow in BERT goes through skip connections instead of attention heads. We further show that consistency of patterns across instances is an indicator of BERT's performance. Finally, we demonstrate that patterns account for far more model performance than previous attention-based and layer-based methods.

[Robust Regression Revisited: Acceleration and Improved Estimation Rates](#)

- Arun Jambulapati, Jerry Li, Tselil Schramm, Kevin Tian
- abstract@[open-review](#): We study fast algorithms for statistical regression problems under the strong contamination model, where the goal is to approximately optimize a generalized linear model (GLM) given adversarially corrupted samples. Prior works in this line of research were based on the robust gradient descent framework of [PrasadSBR20](#), a first-order method using biased gradient queries, or the [Sever](#) framework of [DiakonikolasKK019](#), an iterative outlier-removal method calling a stationary point finder. We present nearly-linear time algorithms for robust regression problems with improved runtime or estimation guarantees compared to the state-of-the-art. For the general case of smooth GLMs (e.g. logistic regression), we show that the robust gradient descent framework of [PrasadSBR20](#) can be accelerated, and show our algorithm extends to optimizing the Moreau envelopes of Lipschitz GLMs (e.g. support vector machines), answering several open questions in the literature. For the well-studied case of robust linear regression, we present an alternative approach obtaining improved estimation rates over prior nearly-linear time algorithms.

Interestingly, our algorithm starts with an identifiability proof introduced in the context of the sum-of-squares algorithm of \cite{BakshiP21}, which achieved optimal error rates while requiring large polynomial runtime and sample complexity. We reinterpret their proof within the Sever framework and obtain a dramatically faster and more sample-efficient algorithm under fewer distributional assumptions.

[Automatic Unsupervised Outlier Model Selection](#)

- Yue Zhao, Ryan Rossi, Leman Akoglu
- abstract@[open-review](#): Given an unsupervised outlier detection task on a new dataset, how can we automatically select a good outlier detection algorithm and its hyperparameter(s) (collectively called a model)? In this work, we tackle the unsupervised outlier model selection (UOMS) problem, and propose MetaOD, a principled, data-driven approach to UOMS based on meta-learning. The UOMS problem is notoriously challenging, as compared to model selection for classification and clustering, since (i) model evaluation is infeasible due to the lack of hold-out data with labels, and (ii) model comparison is infeasible due to the lack of a universal objective function. MetaOD capitalizes on the performances of a large body of detection models on historical outlier detection benchmark datasets, and carries over this prior experience to automatically select an effective model to be employed on a new dataset without any labels, model evaluations or model comparisons. To capture task similarity within our meta-learning framework, we introduce specialized meta-features that quantify outlying characteristics of a dataset. Extensive experiments show that selecting a model by MetaOD significantly outperforms no model selection (e.g. always using the same popular model or the ensemble of many) as well as other meta-learning techniques that we tailored for UOMS. Moreover upon (meta-)training, MetaOD is extremely efficient at test time; selecting from a large pool of 300+ models takes less than 1 second for a new task. We open-source MetaOD and our meta-learning database for practical use and to foster further research on the UOMS problem.

[Pruning Randomly Initialized Neural Networks with Iterative Randomization](#)

- Daiki Chijiwa, Shin'ya Yamaguchi, Yasutoshi Ida, Kenji Umakoshi, Tomohiro INOUE
- abstract@[open-review](#): Pruning the weights of randomly initialized neural networks plays an important role in the context of lottery ticket hypothesis. Ramanujan et al. (2020) empirically showed that only pruning the weights can achieve remarkable performance instead of optimizing the weight values. However, to achieve the same level of performance as the weight optimization, the pruning approach requires more parameters in the networks before pruning and thus more memory space. To overcome this parameter inefficiency, we introduce a novel framework to prune randomly initialized neural networks with iteratively randomizing weight values (IteRand). Theoretically, we prove an approximation theorem in our framework, which indicates that the randomizing operations are provably effective to reduce the required number of the parameters. We also empirically demonstrate the parameter efficiency in multiple experiments on CIFAR-10 and ImageNet.

[Probing Inter-modality: Visual Parsing with Self-Attention for Vision-and-Language Pre-training](#)

- Hongwei Xue, Yupan Huang, Bei Liu, Houwen Peng, Jianlong Fu, Houqiang Li, Jiebo Luo
- abstract@[open-review](#): Vision-Language Pre-training (VLP) aims to learn multi-modal representations from image-text pairs and serves for downstream vision-language tasks in a fine-tuning fashion. The dominant VLP models adopt a CNN-Transformer architecture, which embeds images with a CNN, and then aligns images and text with a Transformer. Visual relationship between visual contents plays an important role in image understanding and is the basic for inter-modal alignment learning. However, CNNs have limitations in visual relation learning due to local receptive field's weakness in modeling long-range dependencies. Thus the two objectives of learning visual relation and inter-modal alignment are encapsulated in the same Transformer network. Such design might restrict the inter-modal alignment learning in the Transformer by ignoring the specialized characteristic of each objective. To tackle this, we propose a fully Transformer visual embedding for VLP to better learn visual relation and further promote inter-modal alignment. Specifically, we propose a metric named Inter-Modality Flow (IMF) to measure the interaction between vision and language modalities (i.e., inter-modality). We also design a novel masking optimization mechanism named Masked Feature Regression (MFR) in Transformer to further promote the inter-modality learning. To the best of our knowledge, this is the first study to explore the benefit of Transformer for visual feature learning in VLP. We verify our method on a wide range of vision-language tasks, including Visual Question Answering (VQA), Visual Entailment and Visual Reasoning. Our approach not only outperforms the state-of-the-art VLP performance, but also shows benefits on the IMF metric.

[Stability and Generalization of Bilevel Programming in Hyperparameter Optimization](#)

- Fan Bao, Guoqiang Wu, Chongxuan LI, Jun Zhu, Bo Zhang
- abstract@[open-review](#): The (gradient-based) bilevel programming framework is widely used in hyperparameter optimization and has achieved excellent performance empirically. Previous theoretical work mainly focuses on its optimization properties, while leaving the analysis on generalization largely open. This paper attempts to address the issue by presenting an expectation bound w.r.t. the validation set based on uniform stability. Our results can explain some mysterious behaviours of the bilevel programming in practice, for instance, overfitting to the validation set. We also present an expectation bound for the classical cross-validation algorithm. Our results suggest that gradient-based algorithms can be better than cross-validation under certain conditions in a theoretical perspective. Furthermore, we prove that regularization terms in both the outer and inner levels can relieve the overfitting problem in gradient-based algorithms. In experiments on feature learning and data reweighting for noisy labels, we corroborate our theoretical findings.

[Regime Switching Bandits](#)

- Xiang Zhou, Yi Xiong, Ningyuan Chen, Xuefeng GAO
- abstract@[open-review](#): We study a multi-armed bandit problem where the rewards exhibit regime switching. Specifically, the distributions of the random rewards generated from all arms are modulated by a common underlying state modeled as a finite-state Markov chain. The agent does not observe the underlying state and has to learn the transition matrix and the reward distributions. We propose a learning algorithm for this problem, building on spectral method-of-moments estimations for hidden Markov models, belief error control in partially observable Markov decision processes and upper-confidence-bound methods for online learning. We also establish an upper bound $\mathcal{O}(T^{2/3}\sqrt{\log T})$ for the proposed learning algorithm where T is the learning horizon. Finally, we conduct proof-of-concept experiments to illustrate the performance of the learning algorithm.

[MixACM: Mixup-Based Robustness Transfer via Distillation of Activated Channel Maps](#)

- Awais Muhammad, Fengwei Zhou, Chuanlong Xie, Jiawei Li, Sung-Ho Bae, Zhenguo Li
- abstract@[open-review](#): Deep neural networks are susceptible to adversarially crafted, small, and imperceptible changes in the natural inputs. The most effective defense mechanism against these examples is adversarial training which constructs adversarial examples during training by iterative maximization of loss. The model is then trained to minimize the loss on these constructed examples. This min-max optimization requires more data, larger capacity models, and additional computing resources. It also degrades the standard generalization performance of a model. Can we achieve robustness more efficiently? In this work, we explore this question from the perspective of knowledge transfer. First, we theoretically show the transferability of robustness from an adversarially trained teacher model to a student model with the help of mixup augmentation. Second, we propose a novel robustness transfer method called Mixup-Based Activated Channel Maps (MixACM) Transfer. MixACM transfers robustness from a robust teacher to a student by matching activated channel maps generated without expensive adversarial perturbations. Finally, extensive experiments on multiple datasets and different learning scenarios show our method can transfer robustness while also improving generalization on natural images.

[Localization, Convexity, and Star Aggregation](#)

- Suhas Vijaykumar
- abstract@[open-review](#): Offset Rademacher complexities have been shown to provide tight upper bounds for the square loss in a broad class of problems including improper statistical learning and online learning. We show that the offset complexity can be generalized to any loss that satisfies a certain general convexity condition. Further, we show that this condition is closely related to both exponential concavity and self-concordance, unifying apparently disparate results. By a novel geometric argument, many of our bounds translate to improper learning in a non-convex class with Audibert's star algorithm. Thus, the offset complexity provides a versatile analytic tool that covers both convex empirical risk minimization and improper learning under entropy conditions. Applying the method, we recover the optimal rates for proper and improper learning with the $\$p\$$ -loss for $1 < p < \infty$, and show that improper variants of empirical risk minimization can attain fast rates for logistic regression and other generalized linear models.

[Aligning Silhouette Topology for Self-Adaptive 3D Human Pose Recovery](#)

- Ramesha Rakesh Mugaludi, Jogendra Nath Kundu, Varun Jampani, Venkatesh Babu R
- abstract@[open-review](#): Articulation-centric 2D/3D pose supervision forms the core training objective in most existing 3D human pose estimation techniques. Except for synthetic source environments, acquiring such rich supervision for each real target domain at deployment is highly inconvenient. However, we realize that standard foreground silhouette estimation techniques (on static camera feeds) remain unaffected by domain-shifts. Motivated by this, we propose a novel target adaptation framework that relies only on silhouette supervision to adapt a source-trained model-based regressor. However, in the absence of any auxiliary cue (multi-view, depth, or 2D pose), an isolated silhouette loss fails to provide a reliable pose-specific gradient and requires to be employed in tandem with a topology-centric loss. To this end, we develop a series of convolution-friendly spatial transformations in order to disentangle a topological-skeleton representation from the raw silhouette. Such a design paves the way to devise a Chamfer-inspired spatial topological-alignment loss via distance field computation, while effectively avoiding any gradient hindering spatial-to-pointset mapping. Experimental results demonstrate our superiority against prior-arts in self-adapting a source trained model to diverse unlabeled target domains, such as a) in-the-wild datasets, b) low-resolution image domains, and c) adversarially perturbed image domains (via UAP).

[Self-Adaptable Point Processes with Nonparametric Time Decays](#)

- Zhimeng Pan, Zheng Wang, Jeff M Phillips, Shandian Zhe
- abstract@[open-review](#): Many applications involve multi-type event data. Understanding the complex influences of the events on each other is critical to discover useful knowledge and to predict future events and their types. Existing methods either ignore or partially account for these influences. Recent works use recurrent neural networks to model the event rate. While being highly expressive, they couple all the temporal dependencies in a black-box and can hardly extract meaningful knowledge. More important, most methods assume an exponential time decay of the influence strength, which is oversimplified and can miss many important strength varying patterns. To overcome these limitations, we propose SPRITE, a $\$S\$$ elf-adaptable $\$P\$$ oint $\$R\$$ ocess w $\$I\$$ th nonparametric $\$T\$$ ime d $\$E\$$ cays, which can decouple the influences between every pair of the events and capture various time decays of the influence strengths. Specifically, we use an embedding to represent each event type and model the event influence as an unknown function of the embeddings and time span. We derive a general construction that can cover all possible time decaying functions. By placing Gaussian process (GP) priors over the latent functions and using Gauss-Legendre quadrature to obtain the integral in the construction, we can flexibly estimate all kinds of time-decaying influences, without restricting to any specific form or imposing derivative constraints that bring learning difficulties. We then use weight space augmentation of GPs to develop an efficient stochastic variational learning algorithm. We show the advantages of our approach in both the ablation study and real-world applications.

[Offline Meta Reinforcement Learning -- Identifiability Challenges and Effective Data Collection Strategies](#)

- Ron Dorfman, Idan Shenfeld, Aviv Tamar
- abstract@[open-review](#): Consider the following instance of the Offline Meta Reinforcement Learning (OMRL) problem: given the complete training logs of $\$N\$$ conventional RL agents, trained on $\$N\$$ different tasks, design a meta-agent that can quickly maximize reward in a new, unseen task from the same task distribution. In particular, while each conventional RL agent explored and exploited its own different task, the meta-agent must identify regularities in the data that lead to effective exploration/exploitation in the unseen task. Here, we take a Bayesian RL (BRL) view, and seek to learn a Bayes-optimal policy from the offline data. Building on the recent VariBAD BRL approach, we develop an off-policy BRL method that learns to plan an exploration strategy based on an adaptive neural belief estimate. However, learning to infer such a belief from offline data brings a new identifiability issue we term MDP ambiguity. We characterize the problem, and suggest resolutions via data collection and modification procedures. Finally, we evaluate our framework on a diverse set of domains, including difficult sparse reward tasks, and demonstrate learning of effective exploration behavior that is qualitatively different from the exploration used by any RL agent in the data. Our code is available online at [url{https://github.com/Rondorf/BOrL}](https://github.com/Rondorf/BOrL).

[RoMA: Robust Model Adaptation for Offline Model-based Optimization](#)

- Sihyun Yu, Sungsoo Ahn, Le Song, Jinwoo Shin
- abstract@[open-review](#): We consider the problem of searching an input maximizing a black-box objective function given a static dataset of input-output queries. A popular approach to solving this problem is maintaining a proxy model, e.g., a deep neural network (DNN), that approximates the true objective function. Here, the main challenge is how to avoid adversarially optimized inputs during the search, i.e., the inputs where the DNN highly overestimates the true objective function. To handle the issue, we propose a new framework, coined robust model adaptation (RoMA), based on gradient-based optimization of inputs over the DNN. Specifically, it consists of two steps: (a) a pre-training strategy to robustly train the proxy model and (b) a novel adaptation procedure of the proxy model to have robust estimates for a specific set of candidate solutions. At a high level, our scheme utilizes the local smoothness prior to overcome the brittleness of the DNN. Experiments under various tasks show the effectiveness of RoMA compared with previous methods, obtaining state-of-the-art results, e.g., RoMA outperforms all at 4 out of 6 tasks and achieves runner-up results at the remaining tasks.

[Flexible Option Learning](#)

- Martin Klissarov, Doina Precup
- abstract@[open-review](#): Temporal abstraction in reinforcement learning (RL), offers the promise of improving generalization and knowledge transfer in complex environments, by propagating information more efficiently over time. Although option learning was initially formulated in a way that allows updating many options simultaneously, using off-policy, intra-option learning (Sutton, Precup & Singh, 1999), many of the recent hierarchical reinforcement learning approaches only update a single option at a time: the option currently executing. We revisit and extend intra-option learning in the context of deep reinforcement learning, in order to enable updating all options consistent with current primitive action choices, without introducing any additional estimates. Our method can therefore be naturally adopted in most hierarchical RL frameworks. When we combine our approach with the option-critic algorithm for option discovery, we obtain significant improvements in performance and data-efficiency across a wide variety of domains.

[Faster Directional Convergence of Linear Neural Networks under Spherically Symmetric Data](#)

- Dachao Lin, Ruoyu Sun, Zhihua Zhang
- abstract@[open-review](#): In this paper, we study gradient methods for training deep linear neural networks with binary cross-entropy loss. In particular, we show global directional convergence guarantees from a polynomial rate to a linear rate for (deep) linear networks with spherically symmetric data distribution, which can be viewed as a specific zero-margin dataset. Our results do not require the assumptions in other works such as small initial loss, presumed convergence of weight direction, or overparameterization. We also characterize our findings in experiments.

[Online Facility Location with Multiple Advice](#)

- Matteo Almanza, Flavio Chierichetti, Silvio Lattanzi, Alessandro Panconesi, Giuseppe Re
- abstract@[open-review](#): Clustering is a central topic in unsupervised learning and its online formulation has received a lot of attention in recent years. In this paper, we study the classic facility location problem in the presence of multiple machine-learned advice. We design an algorithm with provable performance guarantees such that, if the advice is good, it outperforms the best-known online algorithms for the problem, and if it is bad it still matches their performance. We complement our theoretical analysis with an in-depth study of the performance of our algorithm, showing its effectiveness on synthetic and real-world data sets.

[Credit Assignment in Neural Networks through Deep Feedback Control](#)

- Alexander Meulemans, Matilde Tristany Farinha, Javier Garcia Ordóñez, Pau Vilimelis Aceituno, João Sacramento, Benjamin F. Grewe
- abstract@[open-review](#): The success of deep learning sparked interest in whether the brain learns by using similar techniques for assigning credit to each synaptic weight for its contribution to the network output. However, the majority of current attempts at biologically-plausible learning methods are either non-local in time, require highly specific connectivity motifs, or have no clear link to any known mathematical optimization method. Here, we introduce Deep Feedback Control (DFC), a new learning method that uses a feedback controller to drive a deep neural network to match a desired output target and whose control signal can be used for credit assignment. The resulting learning rule is fully local in space and time and approximates Gauss-Newton optimization for a wide range of feedback connectivity patterns. To further underline its biological plausibility, we relate DFC to a multi-compartment model of cortical pyramidal neurons with a local voltage-dependent synaptic plasticity rule, consistent with recent theories of dendritic processing. By combining dynamical system theory with mathematical optimization theory, we provide a strong theoretical foundation for DFC that we corroborate with detailed results on toy experiments and standard computer-vision benchmarks.

[Robust Online Correlation Clustering](#)

- Silvio Lattanzi, Benjamin Moseley, Sergei Vassilvitskii, Yuyan Wang, Rudy Zhou
- abstract@[open-review](#): In correlation clustering we are given a set of points along with recommendations whether each pair of points should be placed in the same cluster or into separate clusters. The goal cluster the points to minimize disagreements from the recommendations. We study the correlation clustering problem in the online setting., where points arrive one at a time, and upon arrival the algorithm must make an irrevocable cluster assignment decision. While the online version is natural, there is a simple lower bound that rules out any algorithm with a non-trivial competitive ratio. In this work we go beyond worst case analysis, and show that the celebrated Pivot algorithm performs well when given access to a small number of random samples from the input. Moreover, we prove that Pivot is robust to additional adversarial perturbations of the sample set in this setting. We conclude with an empirical analysis validating our theoretical findings.

[Neural Additive Models: Interpretable Machine Learning with Neural Nets](#)

- Rishabh Agarwal, Levi Melnick, Nicholas Frosst, Xuezhou Zhang, Ben Lengerich, Rich Caruana, Geoffrey E. Hinton
- abstract@[open-review](#): Deep neural networks (DNNs) are powerful black-box predictors that have achieved impressive performance on a wide variety of tasks. However, their accuracy comes at the cost of intelligibility: it is usually unclear how they make their decisions. This hinders their applicability to high stakes decision-making domains such as healthcare. We propose Neural Additive Models (NAMs) which combine some of the expressivity of DNNs with the inherent intelligibility of generalized additive models. NAMs learn a linear combination of neural networks that each attend to a single input feature. These networks are trained jointly and can learn arbitrarily complex relationships between their input feature and the output. Our experiments on regression and classification datasets show that NAMs are more accurate than widely used intelligible models such as logistic regression and shallow decision trees. They perform similarly to existing state-of-the-art generalized additive models in accuracy, but are more flexible because they are based on neural nets instead of boosted trees. To demonstrate this, we show how NAMs can be used for multitask learning on synthetic data and on the COMPAS recidivism data due to their composability, and demonstrate that the differentiability of NAMs allows them to train more complex interpretable models for COVID-19.

[Representation Learning for Event-based Visuomotor Policies](#)

- Sai Venkrala, Sami Mian, Ashish Kapoor
- abstract@[open-review](#): Event-based cameras are dynamic vision sensors that provide asynchronous measurements of changes in per-pixel brightness at a microsecond level. This makes them significantly faster than conventional frame-based cameras, and an appealing choice for high-speed robot navigation. While an interesting sensor modality, this asynchronously streamed event data poses a challenge for machine learning based computer vision techniques that are more suited for synchronous, frame-based data. In this paper, we present an event variational autoencoder through which compact representations can be learnt directly from asynchronous spatiotemporal event data. Furthermore, we show that such pretrained representations can be used for event-based reinforcement learning instead of end-to-end reward driven perception. We validate this framework of learning event-based visuomotor policies by applying it to an obstacle avoidance scenario in simulation. Compared to techniques that treat event data as images, we show that representations learnt from event streams result in faster policy training, adapt to different control capacities, and demonstrate a higher degree of robustness to environmental changes and sensor noise.

[Kernel Functional Optimisation](#)

- Arun Kumar Anjanapura Venkatesh, Alistair Shilton, Santu Rana, Sunil Gupta, Svetha Venkatesh
- abstract@[open-review](#): Traditional methods for kernel selection rely on parametric kernel functions or a combination thereof and although the kernel hyperparameters are tuned, these methods often provide sub-optimal results due to the limitations induced by the parametric forms. In this paper, we propose a novel formulation for kernel selection using efficient Bayesian optimisation to find the best fitting non-parametric kernel. The kernel is expressed using a linear combination of functions sampled from a prior Gaussian Process (GP) defined by a hyperkernel. We also provide a mechanism to ensure the positive definiteness of the Gram matrix constructed using the resultant kernels. Our experimental results on GP regression and Support Vector Machine (SVM) classification tasks involving both synthetic functions and several real-world datasets show the superiority of our approach over the state-of-the-art.

[Generalized Shape Metrics on Neural Representations](#)

- Alex H Williams, Erin Kunz, Simon Kornblith, Scott Linderman
- abstract@[open-review](#): Understanding the operation of biological and artificial networks remains a difficult and important challenge. To identify general principles, researchers are increasingly interested in surveying large collections of networks that are trained on, or biologically adapted to, similar tasks. A standardized set of analysis tools is now needed to identify how network-level covariates--such as architecture, anatomical brain region, and model organism--impact neural representations (hidden layer activations). Here, we provide a rigorous foundation for these analyses by defining a broad family of metric spaces that quantify representational dissimilarity. Using this framework, we modify existing representational similarity measures based on canonical correlation analysis and centered kernel alignment to satisfy the triangle inequality, formulate a novel metric that respects the inductive biases in convolutional layers, and identify approximate Euclidean embeddings that enable network representations to be incorporated into essentially any off-the-

shelf machine learning method. We demonstrate these methods on large-scale datasets from biology (Allen Institute Brain Observatory) and deep learning (NAS-Bench-101). In doing so, we identify relationships between neural representations that are interpretable in terms of anatomical features and model performance.

[Diverse Message Passing for Attribute with Heterophily](#)

- Liang Yang, Mengzhe Li, Liyang Liu, bingxin niu, Chuan Wang, Xiaochun Cao, Yuanfang Guo
- abstract@[open-review](#): Most of the existing GNNs can be modeled via the Uniform Message Passing framework. This framework considers all the attributes of each node in its entirety, shares the uniform propagation weights along each edge, and focuses on the uniform weight learning. The design of this framework possesses two prerequisites, the simplification of homophily and heterophily to the node-level property and the ignorance of attribute differences. Unfortunately, different attributes possess diverse characteristics. In this paper, the network homophily rate defined with respect to the node labels is extended to attribute homophily rate by taking the attributes as weak labels. Based on this attribute homophily rate, we propose a Diverse Message Passing (DMP) framework, which specifies every attribute propagation weight on each edge. Besides, we propose two specific strategies to significantly reduce the computational complexity of DMP to prevent the overfitting issue. By investigating the spectral characteristics, existing spectral GNNs are actually equivalent to a degenerated version of DMP. From the perspective of numerical optimization, we provide a theoretical analysis to demonstrate DMP's powerful representation ability and the ability of alleviating the over-smoothing issue. Evaluations on various real networks demonstrate the superiority of our DMP on handling the networks with heterophily and alleviating the over-smoothing issue, compared to the existing state-of-the-arts.

[Towards Robust Bisimulation Metric Learning](#)

- Mete Kemertas, Tristan Aumentado-Armstrong
- abstract@[open-review](#): Learned representations in deep reinforcement learning (DRL) have to extract task-relevant information from complex observations, balancing between robustness to distraction and informativeness to the policy. Such stable and rich representations, often learned via modern function approximation techniques, can enable practical application of the policy improvement theorem, even in high-dimensional continuous state-action spaces. Bisimulation metrics offer one solution to this representation learning problem, by collapsing functionally similar states together in representation space, which promotes invariance to noise and distractors. In this work, we generalize value function approximation bounds for on-policy bisimulation metrics to non-optimal policies and approximate environment dynamics. Our theoretical results help us identify embedding pathologies that may occur in practical use. In particular, we find that these issues stem from an underconstrained dynamics model and an unstable dependence of the embedding norm on the reward signal in environments with sparse rewards. Further, we propose a set of practical remedies: (i) a norm constraint on the representation space, and (ii) an extension of prior approaches with intrinsic rewards and latent space regularization. Finally, we provide evidence that the resulting method is not only more robust to sparse reward functions, but also able to solve challenging continuous control tasks with observational distractions, where prior methods fail.

[Beyond BatchNorm: Towards a Unified Understanding of Normalization in Deep Learning](#)

- Ekdeep S Lubana, Robert Dick, Hidenori Tanaka
- abstract@[open-review](#): Inspired by BatchNorm, there has been an explosion of normalization layers in deep learning. Recent works have identified a multitude of beneficial properties in BatchNorm to explain its success. However, given the pursuit of alternative normalization layers, these properties need to be generalized so that any given layer's success/failure can be accurately predicted. In this work, we take a first step towards this goal by extending known properties of BatchNorm in randomly initialized deep neural networks (DNNs) to several recently proposed normalization layers. Our primary findings follow: (i) similar to BatchNorm, activations-based normalization layers can prevent exponential growth of activations in ResNets, but parametric techniques require explicit remedies; (ii) use of GroupNorm can ensure an informative forward propagation, with different samples being assigned dissimilar activations, but increasing group size results in increasingly indistinguishable activations for different samples, explaining slow convergence speed in models with LayerNorm; and (iii) small group sizes result in large gradient norm in earlier layers, hence explaining training instability issues in Instance Normalization and illustrating a speed-stability tradeoff in GroupNorm. Overall, our analysis reveals a unified set of mechanisms that underpin the success of normalization methods in deep learning, providing us with a compass to systematically explore the vast design space of DNN normalization layers.

[Representation Learning Beyond Linear Prediction Functions](#)

- Ziping Xu, Ambuj Tewari
- abstract@[open-review](#): Recent papers on the theory of representation learning has shown the importance of a quantity called diversity when generalizing from a set of source tasks to a target task. Most of these papers assume that the function mapping shared representations to predictions is linear, for both source and target tasks. In practice, researchers in deep learning use different numbers of extra layers following the pretrained model based on the difficulty of the new task. This motivates us to ask whether diversity can be achieved when source tasks and the target task use different prediction function spaces beyond linear functions. We show that diversity holds even if the target task uses a neural network with multiple layers, as long as source tasks use linear functions. If source tasks use nonlinear prediction functions, we provide a negative result by showing that depth-1 neural networks with ReLu activation function need exponentially many source tasks to achieve diversity. For a general function class, we find that eluder dimension gives a lower bound on the number of tasks required for diversity. Our theoretical results imply that simpler tasks generalize better. Though our theoretical results are shown for the global minimizer of empirical risks, their qualitative predictions still hold true for gradient-based optimization algorithms as verified by our simulations on deep neural networks.

[Volume Rendering of Neural Implicit Surfaces](#)

- Lior Yariv, Jiatao Gu, Yoni Kasten, Yaron Lipman
- abstract@[open-review](#): Neural volume rendering became increasingly popular recently due to its success in synthesizing novel views of a scene from a sparse set of input images. So far, the geometry learned by neural volume rendering techniques was modeled using a generic density function. Furthermore, the geometry itself was extracted using an arbitrary level set of the density function leading to a noisy, often low fidelity reconstruction. The goal of this paper is to improve geometry representation and reconstruction in neural volume rendering. We achieve that by modeling the volume density as a function of the geometry. This is in contrast to previous work modeling the geometry as a function of the volume density. In more detail, we define the volume density function as Laplace's cumulative distribution function (CDF) applied to a signed distance function (SDF) representation. This simple density representation has three benefits: (i) it provides a useful inductive bias to the geometry learned in the neural volume rendering process; (ii) it facilitates a bound on the opacity approximation error, leading to an accurate sampling of the viewing ray. Accurate sampling is important to provide a precise coupling of geometry and radiance; and (iii) it allows efficient unsupervised disentanglement of shape and appearance in volume rendering. Applying this new density representation to challenging scene multiview datasets produced high quality geometry reconstructions, outperforming relevant baselines. Furthermore, switching shape and appearance between scenes is possible due to the disentanglement of the two.

[MAUVE: Measuring the Gap Between Neural Text and Human Text using Divergence Frontiers](#)

- Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, Zaid Harchaoui

- abstract@[open-review](#): As major progress is made in open-ended text generation, measuring how close machine-generated text is to human language remains a critical open problem. We introduce Mauve, a comparison measure for open-ended text generation, which directly compares the learnt distribution from a text generation model to the distribution of human-written text using divergence frontiers. Mauve scales up to modern text generation models by computing information divergences in a quantized embedding space. Through an extensive empirical study on three open-ended generation tasks, we find that Mauve identifies known properties of generated text, scales naturally with model size, and correlates with human judgments, with fewer restrictions than existing distributional evaluation metrics.

[Accurately Solving Rod Dynamics with Graph Learning](#)

- Han Shao, Tassilo Kugelstadt, Torsten Hädrich, Wojtek Palubicki, Jan Bender, Soeren Pirk, Dominik L Michels
- abstract@[open-review](#): Iterative solvers are widely used to accurately simulate physical systems. These solvers require initial guesses to generate a sequence of improving approximate solutions. In this contribution, we introduce a novel method to accelerate iterative solvers for rod dynamics with graph networks (GNs) by predicting the initial guesses to reduce the number of iterations. Unlike existing methods that aim to learn physical systems in an end-to-end manner, our approach guarantees long-term stability and therefore leads to more accurate solutions. Furthermore, our method improves the run time performance of traditional iterative solvers for rod dynamics. To explore our method we make use of position-based dynamics (PBD) as a common solver for physical systems and evaluate it by simulating the dynamics of elastic rods. Our approach is able to generalize across different initial conditions, discretizations, and realistic material properties. We demonstrate that it also performs well when taking discontinuous effects into account such as collisions between individual rods. Finally, to illustrate the scalability of our approach, we simulate complex 3D tree models composed of over a thousand individual branch segments swaying in wind fields.

[Limiting fluctuation and trajectory stability of multilayer neural networks with mean field training](#)

- Huy Tuan Pham, Phan-Minh Nguyen
- abstract@[open-review](#): The mean field theory of multilayer neural networks centers around a particular infinite-width scaling, in which the learning dynamics is shown to be closely tracked by the mean field limit. A random fluctuation around this infinite-width limit is expected from a large-width expansion to the next order. This fluctuation has been studied only in the case of shallow networks, where previous works employ heavily technical notions or additional formulation ideas amenable only to that case. Treatment of the multilayer case has been missing, with the chief difficulty in finding a formulation that must capture the stochastic dependency across not only time but also depth. In this work, we initiate the study of the fluctuation in the case of multilayer networks, at any network depth. Leveraging on the neuronal embedding framework recently introduced by Nguyen and Pham, we systematically derive a system of dynamical equations, called the second-order mean field limit, that captures the limiting fluctuation distribution. We demonstrate through the framework the complex interaction among neurons in this second-order mean field limit, the stochasticity with cross-layer dependency and the nonlinear time evolution inherent in the limiting fluctuation. A limit theorem is proven to relate quantitatively this limit to the fluctuation realized by large-width networks. We apply the result to show a stability property of gradient descent mean field training: in the large-width regime, along the training trajectory, it progressively biases towards a solution with "minimal fluctuation" (in fact, vanishing fluctuation) in the learned output function, even after the network has been initialized at or has converged (sufficiently fast) to a global optimum. This extends a similar phenomenon previously shown only for shallow networks with a squared loss in the empirical risk minimization setting, to multilayer networks with a loss function that is not necessarily convex in a more general setting.

[Medical Dead-ends and Learning to Identify High-Risk States and Treatments](#)

- Mehdi Fatemi, Taylor W. Killian, Jayakumar Subramanian, Marzyeh Ghassemi
- abstract@[open-review](#): Machine learning has successfully framed many sequential decision making problems as either supervised prediction, or optimal decision-making policy identification via reinforcement learning. In data-constrained offline settings, both approaches may fail as they assume fully optimal behavior or rely on exploring alternatives that may not exist. We introduce an inherently different approach that identifies "dead-ends" of a state space. We focus on patient condition in the intensive care unit, where a "medical dead-end" indicates that a patient will expire, regardless of all potential future treatment sequences. We postulate "treatment security" as avoiding treatments with probability proportional to their chance of leading to dead-ends, present a formal proof, and frame discovery as an RL problem. We then train three independent deep neural models for automated state construction, dead-end discovery and confirmation. Our empirical results discover that dead-ends exist in real clinical data among septic patients, and further reveal gaps between secure treatments and those administered.

[Overcoming the Convex Barrier for Simplex Inputs](#)

- Harkirat Singh Behl, M. Pawan Kumar, Philip Torr, Krishnamurthy Dvijotham
- abstract@[open-review](#): Recent progress in neural network verification has challenged the notion of a convex barrier, that is, an inherent weakness in the convex relaxation of the output of a neural network. Specifically, there now exists a tight relaxation for verifying the robustness of a neural network to ℓ_∞ input perturbations, as well as efficient primal and dual solvers for the relaxation. Buoyed by this success, we consider the problem of developing similar techniques for verifying robustness to input perturbations within the probability simplex. We prove a somewhat surprising result that, in this case, not only can one design a tight relaxation that overcomes the convex barrier, but the size of the relaxation remains linear in the number of neurons, thereby leading to simpler and more efficient algorithms. We establish the scalability of our overall approach via the specification of ℓ_1 robustness for CIFAR-10 and MNIST classification, where our approach improves the state of the art verified accuracy by up to 14.4%. Furthermore, we establish its accuracy on a novel and highly challenging task of verifying the robustness of a multi-modal (text and image) classifier to arbitrary changes in its textual input.

[High-probability Bounds for Non-Convex Stochastic Optimization with Heavy Tails](#)

- Ashok Cutkosky, Harsh Mehta
- abstract@[open-review](#): We consider non-convex stochastic optimization using first-order algorithms for which the gradient estimates may have heavy tails. We show that a combination of gradient clipping, momentum, and normalized gradient descent yields convergence to critical points in high-probability with best-known rates for smooth losses when the gradients only have bounded $\mathbb{E}[\|\mathbf{g}\|_p^p]$ moments for some $p \in (1, 2)$. We then consider the case of second-order smooth losses, which to our knowledge have not been studied in this setting, and again obtain high-probability bounds for any $\mathbb{E}[\|\mathbf{g}\|_p^p]$. Moreover, our results hold for arbitrary smooth norms, in contrast to the typical SGD analysis which requires a Hilbert space norm. Further, we show that after a suitable "burn-in" period, the objective value will monotonically decrease for every iteration until a critical point is identified, which provides intuition behind the popular practice of learning rate "warm-up" and also yields a last-iterate guarantee.

[Batch Normalization Orthogonalizes Representations in Deep Random Networks](#)

- Hadi Daneshmand, Amir Joudaki, Francis Bach
- abstract@[open-review](#): This paper underlines an elegant property of batch-normalization (BN): Successive batch normalizations with random linear updates make samples increasingly orthogonal. We establish a non-asymptotic characterization of the interplay between depth, width, and the orthogonality of deep representations. More precisely, we prove, under a mild assumption, the deviation of the representations from orthogonality rapidly decays with depth up to a term inversely proportional to the network width. This result has two main theoretical and practical implications: 1) Theoretically, as the depth grows, the distribution of the outputs contracts to a Wasserstein-2 ball around an isotropic normal distribution. Furthermore, the

radius of this Wasserstein ball shrinks with the width of the network. 2) Practically, the orthogonality of the representations directly influences the performance of stochastic gradient descent (SGD). When representations are initially aligned, we observe SGD wastes many iterations to disentangle representations before the classification. Nevertheless, we experimentally show that starting optimization from orthogonal representations is sufficient to accelerate SGD, with no need for BN.

[Support vector machines and linear regression coincide with very high-dimensional features](#)

- Navid Ardestir, Clayton Sanford, Daniel J. Hsu
- abstract@[open-review](#): The support vector machine (SVM) and minimum Euclidean norm least squares regression are two fundamentally different approaches to fitting linear models, but they have recently been connected in models for very high-dimensional data through a phenomenon of support vector proliferation, where every training example used to fit an SVM becomes a support vector. In this paper, we explore the generality of this phenomenon and make the following contributions. First, we prove a super-linear lower bound on the dimension (in terms of sample size) required for support vector proliferation in independent feature models, matching the upper bounds from previous works. We further identify a sharp phase transition in Gaussian feature models, bound the width of this transition, and give experimental support for its universality. Finally, we hypothesize that this phase transition occurs only in much higher-dimensional settings in the $\|\cdot\|_1$ variant of the SVM, and we present a new geometric characterization of the problem that may elucidate this phenomenon for the general $\|\cdot\|_p$ case.

[Coupled Segmentation and Edge Learning via Dynamic Graph Propagation](#)

- Zhiding Yu, Rui Huang, Wonmin Byeon, Sifei Liu, Guilin Liu, Thomas Breuel, Anima Anandkumar, Jan Kautz
- abstract@[open-review](#): Image segmentation and edge detection are both central problems in perceptual grouping. It is therefore interesting to study how these two tasks can be coupled to benefit each other. Indeed, segmentation can be easily transformed into contour edges to guide edge learning. However, the converse is nontrivial since general edges may not always form closed contours. In this paper, we propose a principled end-to-end framework for coupled edge and segmentation learning, where edges are leveraged as pairwise similarity cues to guide segmentation. At the core of our framework is a recurrent module termed as dynamic graph propagation (DGP) layer that performs message passing on dynamically constructed graphs. The layer uses learned gating to dynamically select neighbors for message passing using max-pooling. The output from message passing is further gated with an edge signal to refine segmentation. Experiments demonstrate that the proposed framework is able to let both tasks mutually improve each other. On Cityscapes validation, our best model achieves 83.7% mIoU in semantic segmentation and 78.7% maximum F-score in semantic edge detection. Our method also leads to improved zero-shot robustness on Cityscapes with natural corruptions (Cityscapes-C).

[Offline RL Without Off-Policy Evaluation](#)

- David Brandfonbrener, Will Whitney, Rajesh Ranganath, Joan Bruna
- abstract@[open-review](#): Most prior approaches to offline reinforcement learning (RL) have taken an iterative actor-critic approach involving off-policy evaluation. In this paper we show that simply doing one step of constrained/regularized policy improvement using an on-policy Q estimate of the behavior policy performs surprisingly well. This one-step algorithm beats the previously reported results of iterative algorithms on a large portion of the D4RL benchmark. The one-step baseline achieves this strong performance while being notably simpler and more robust to hyperparameters than previously proposed iterative algorithms. We argue that the relatively poor performance of iterative approaches is a result of the high variance inherent in doing off-policy evaluation and magnified by the repeated optimization of policies against those estimates. In addition, we hypothesize that the strong performance of the one-step algorithm is due to a combination of favorable structure in the environment and behavior policy.

[Continuous vs. Discrete Optimization of Deep Neural Networks](#)

- Omer Elkabetz, Nadav Cohen
- abstract@[open-review](#): Existing analyses of optimization in deep learning are either continuous, focusing on (variants of) gradient flow, or discrete, directly treating (variants of) gradient descent. Gradient flow is amenable to theoretical analysis, but is stylized and disregards computational efficiency. The extent to which it represents gradient descent is an open question in the theory of deep learning. The current paper studies this question. Viewing gradient descent as an approximate numerical solution to the initial value problem of gradient flow, we find that the degree of approximation depends on the curvature around the gradient flow trajectory. We then show that over deep neural networks with homogeneous activations, gradient flow trajectories enjoy favorable curvature, suggesting they are well approximated by gradient descent. This finding allows us to translate an analysis of gradient flow over deep linear neural networks into a guarantee that gradient descent efficiently converges to global minimum almost surely under random initialization. Experiments suggest that over simple deep neural networks, gradient descent with conventional step size is indeed close to gradient flow. We hypothesize that the theory of gradient flows will unravel mysteries behind deep learning.

[CrypTen: Secure Multi-Party Computation Meets Machine Learning](#)

- Brian Knott, Shobha Venkataraman, Awni Hannun, Shubho Sengupta, Mark Ibrahim, Laurens van der Maaten
- abstract@[open-review](#): Secure multi-party computation (MPC) allows parties to perform computations on data while keeping that data private. This capability has great potential for machine-learning applications: it facilitates training of machine-learning models on private data sets owned by different parties, evaluation of one party's private model using another party's private data, etc. Although a range of studies implement machine-learning models via secure MPC, such implementations are not yet mainstream. Adoption of secure MPC is hampered by the absence of flexible software frameworks that "speak the language" of machine-learning researchers and engineers. To foster adoption of secure MPC in machine learning, we present CrypTen: a software framework that exposes popular secure MPC primitives via abstractions that are common in modern machine-learning frameworks, such as tensor computations, automatic differentiation, and modular neural networks. This paper describes the design of CrypTen and measure its performance on state-of-the-art models for text classification, speech recognition, and image classification. Our benchmarks show that CrypTen's GPU support and high-performance communication between (an arbitrary number of) parties allows it to perform efficient private evaluation of modern machine-learning models under a semi-honest threat model. For example, two parties using CrypTen can securely predict phonemes in speech recordings using Wav2Letter faster than real-time. We hope that CrypTen will spur adoption of secure MPC in the machine-learning community.

[Can contrastive learning avoid shortcut solutions?](#)

- Joshua Robinson, Li Sun, Ke Yu, Kayhan Batmanghelich, Stefanie Jegelka, Suvrit Sra
- abstract@[open-review](#): The generalization of representations learned via contrastive learning depends crucially on what features of the data are extracted. However, we observe that the contrastive loss does not always sufficiently guide which features are extracted, a behavior that can negatively impact the performance on downstream tasks via "shortcuts", i.e., by inadvertently suppressing important predictive features. We find that feature extraction is influenced by the difficulty of the so-called instance discrimination task (i.e., the task of discriminating pairs of similar points from pairs of dissimilar ones). Although harder pairs improve the representation of some features, the improvement comes at the cost of suppressing previously well represented features. In response, we propose implicit feature modification (IFM), a method for altering positive and negative samples in order to guide contrastive models towards capturing a wider variety of predictive features. Empirically, we observe that IFM reduces feature suppression, and as a result improves performance on vision and medical imaging tasks.

[See More for Scene: Pairwise Consistency Learning for Scene Classification](#)

- Gongwei Chen, Xinhang Song, Bohan Wang, Shuqiang Jiang
- abstract@[open-review](#): Scene classification is a valuable classification subtask and has its own characteristics which still needs more in-depth studies. Basically, scene characteristics are distributed over the whole image, which cause the need of “seeing” comprehensive and informative regions. Previous works mainly focus on region discovery and aggregation, while rarely involves the inherent properties of CNN along with its potential ability to satisfy the requirements of scene classification. In this paper, we propose to understand scene images and the scene classification CNN models in terms of the focus area. From this new perspective, we find that large focus area is preferred in scene classification CNN models as a consequence of learning scene characteristics. Meanwhile, the analysis about existing training schemes helps us to understand the effects of focus area, and also raises the question about optimal training method for scene classification. Pursuing the better usage of scene characteristics, we propose a new learning scheme with a tailored loss in the goal of activating larger focus area on scene images. Since the supervision of the target regions to be enlarged is usually lacked, our alternative learning scheme is to erase already activated area, and allow the CNN models to activate more area during training. The proposed scheme is implemented by keeping the pairwise consistency between the output of the erased image and its original one. In particular, a tailored loss is proposed to keep such pairwise consistency by leveraging category-relevance information. Experiments on Places365 show the significant improvements of our method with various CNNs. Our method shows an inferior result on the object-centric dataset, ImageNet, which experimentally indicates that it captures the unique characteristics of scenes.

[Provable Guarantees for Self-Supervised Deep Learning with Spectral Contrastive Loss](#)

- Jeff Z. HaoChen, Colin Wei, Adrien Gaidon, Tengyu Ma
- abstract@[open-review](#): Recent works in self-supervised learning have advanced the state-of-the-art by relying on the contrastive learning paradigm, which learns representations by pushing positive pairs, or similar examples from the same class, closer together while keeping negative pairs far apart. Despite the empirical successes, theoretical foundations are limited -- prior analyses assume conditional independence of the positive pairs given the same class label, but recent empirical applications use heavily correlated positive pairs (i.e., data augmentations of the same image). Our work analyzes contrastive learning without assuming conditional independence of positive pairs using a novel concept of the augmentation graph on data. Edges in this graph connect augmentations of the same data, and ground-truth classes naturally form connected sub-graphs. We propose a loss that performs spectral decomposition on the population augmentation graph and can be succinctly written as a contrastive learning objective on neural net representations. Minimizing this objective leads to features with provable accuracy guarantees under linear probe evaluation. By standard generalization bounds, these accuracy guarantees also hold when minimizing the training contrastive loss. In all, this work provides the first provable analysis for contrastive learning where the guarantees can apply to realistic empirical settings.

[Greedy Approximation Algorithms for Active Sequential Hypothesis Testing](#)

- Kyra Gan, Su Jia, Andrew Li
- abstract@[open-review](#): In the problem of \emph{active sequential hypothesis testing} (ASHT), a learner seeks to identify the \emph{true} hypothesis from among a known set of hypotheses. The learner is given a set of actions and knows the random distribution of the outcome of any action under any true hypothesis. Given a target error $\$\\delta>0$, the goal is to sequentially select the fewest number of actions so as to identify the true hypothesis with probability at least $1 - \\delta$. Motivated by applications in which the number of hypotheses or actions is massive (e.g., genomics-based cancer detection), we propose efficient (greedy, in fact) algorithms and provide the first approximation guarantees for ASHT, under two types of adaptivity. Both of our guarantees are independent of the number of actions and logarithmic in the number of hypotheses. We numerically evaluate the performance of our algorithms using both synthetic and real-world DNA mutation data, demonstrating that our algorithms outperform previously proposed heuristic policies by large margins.$$

[When False Positive is Intolerant: End-to-End Optimization with Low FPR for Multipartite Ranking](#)

- Peisong Wen, Qianqian Xu, Zhiyong Yang, Yuan He, Qingming Huang
- abstract@[open-review](#): Multipartite ranking is a basic task in machine learning, where the Area Under the receiver operating characteristics Curve (AUC) is generally applied as the evaluation metric. Despite that AUC reflects the overall performance of the model, it is inconsistent with the expected performance in some application scenarios, where only a low False Positive Rate (FPR) is meaningful. To leverage high performance under low FPRs, we consider an alternative metric for multipartite ranking evaluating the True Positive Rate (TPR) at a given FPR, denoted as TPR@FPR. Unfortunately, the key challenge of direct TPR@FPR optimization is two-fold: \textbf{(a)} the original objective function is not differentiable, making gradient backpropagation impossible; \textbf{(b)} the loss function could not be written as a sum of independent instance-wise terms, making mini-batch based optimization infeasible. To address these issues, we propose a novel framework on top of the deep learning framework named \textit{Cross-Batch Approximation for Multipartite Ranking (CBA-MR)}. In face of \textbf{(a)}, we propose a differentiable surrogate optimization problem where the instances having a short-time effect on FPR are rendered with different weights based on the random walk hypothesis. To tackle \textbf{(b)}, we propose a fast ranking estimation method, where the full-batch loss evaluation is replaced by a delayed update scheme with the help of an embedding cache. Finally, experimental results on four real-world benchmarks are provided to demonstrate the effectiveness of the proposed method.

[Convex Polytope Trees](#)

- Mohammadreza Armandpour, Ali Sadeghian, Mingyuan Zhou
- abstract@[open-review](#): A decision tree is commonly restricted to use a single hyperplane to split the covariate space at each of its internal nodes. It often requires a large number of nodes to achieve high accuracy. In this paper, we propose convex polytope trees (CPT) to expand the family of decision trees by an interpretable generalization of their decision boundary. The splitting function at each node of CPT is based on the logical disjunction of a community of differently weighted probabilistic linear decision-makers, which also geometrically corresponds to a convex polytope in the covariate space. We use a nonparametric Bayesian prior at each node to infer the community's size, encouraging simpler decision boundaries by shrinking the number of polytope facets. We develop a greedy method to efficiently construct CPT and scalable end-to-end training algorithms for the tree parameters when the tree structure is given. We empirically demonstrate the efficiency of CPT over existing state-of-the-art decision trees in several real-world classification and regression tasks from diverse domains.

[The Skellam Mechanism for Differentially Private Federated Learning](#)

- Naman Agarwal, Peter Kairouz, Ziyu Liu
- abstract@[open-review](#): We introduce the multi-dimensional Skellam mechanism, a discrete differential privacy mechanism based on the difference of two independent Poisson random variables. To quantify its privacy guarantees, we analyze the privacy loss distribution via a numerical evaluation and provide a sharp bound on the Rényi divergence between two shifted Skellam distributions. While useful in both centralized and distributed privacy applications, we investigate how it can be applied in the context of federated learning with secure aggregation under communication constraints. Our theoretical findings and extensive experimental evaluations demonstrate that the Skellam mechanism provides the same privacy-accuracy trade-offs as the continuous Gaussian mechanism, even when the precision is low. More importantly, Skellam is closed under summation and sampling from it only requires sampling from a Poisson distribution -- an efficient routine that ships with all machine learning and data analysis software packages. These features, along with its discrete nature and competitive privacy-accuracy trade-offs, make it an attractive practical alternative to the newly introduced discrete Gaussian mechanism.

Stability and Deviation Optimal Risk Bounds with Convergence Rate $\$O(1/n)$

- Yegor Klochkov, Nikita Zhivotovskiy
- abstract@[open-review](#): The sharpest known high probability generalization bounds for uniformly stable algorithms (Feldman, Vondrak, NeurIPS 2018, COLT, 2019), (Bousquet, Klochkov, Zhivotovskiy, COLT, 2020) contain a generally inevitable sampling error term of order $\$\\Theta(1/sqrt{n})$. When applied to excess risk bounds, this leads to suboptimal results in several standard stochastic convex optimization problems. We show that if the so-called Bernstein condition is satisfied, the term $\$\\Theta(1/sqrt{n})$ can be avoided, and high probability excess risk bounds of order up to $\$O(1/n)$ are possible via uniform stability. Using this result, we show a high probability excess risk bound with the rate $\$O(log n/n)$ for strongly convex and Lipschitz losses valid for \emph{any} empirical risk minimization method. This resolves a question of Shalev-Shwartz, Shamir, Srebro, and Sridharan (COLT, 2009). We discuss how $\$O(log n/n)$ high probability excess risk bounds are possible for projected gradient descent in the case of strongly convex and Lipschitz losses without the usual smoothness assumption.

SketchGen: Generating Constrained CAD Sketches

- Wamiq Para, Shariq Bhat, Paul Guerrero, Tom Kelly, Niloy Mitra, Leonidas J. Guibas, Peter Wonka
- abstract@[open-review](#): Computer-aided design (CAD) is the most widely used modeling approach for technical design. The typical starting point in these designs is 2D sketches which can later be extruded and combined to obtain complex three-dimensional assemblies. Such sketches are typically composed of parametric primitives, such as points, lines, and circular arcs, augmented with geometric constraints linking the primitives, such as coincidence, parallelism, or orthogonality. Sketches can be represented as graphs, with the primitives as nodes and the constraints as edges. Training a model to automatically generate CAD sketches can enable several novel workflows, but is challenging due to the complexity of the graphs and the heterogeneity of the primitives and constraints. In particular, each type of primitive and constraint may require a record of different size and parameter types. We propose SketchGen as a generative model based on a transformer architecture to address the heterogeneity problem by carefully designing a sequential language for the primitives and constraints that allows distinguishing between different primitive or constraint types and their parameters, while encouraging our model to re-use information across related parameters, encoding shared structure. A particular highlight of our work is the ability to produce primitives linked via constraints that enables the final output to be further regularized via a constraint solver. We evaluate our model by demonstrating constraint prediction for given sets of primitives and full sketch generation from scratch, showing that our approach significantly outperforms the state-of-the-art in CAD sketch generation.

CLDA: Contrastive Learning for Semi-Supervised Domain Adaptation

- Ankit Singh
- abstract@[open-review](#): Unsupervised Domain Adaptation (UDA) aims to align the labeled source distribution with the unlabeled target distribution to obtain domain invariant predictive models. However, the application of well-known UDA approaches does not generalize well in Semi-Supervised Domain Adaptation (SSDA) scenarios where few labeled samples from the target domain are available. This paper proposes a simple Contrastive Learning framework for semi-supervised Domain Adaptation (CLDA) that attempts to bridge the intra-domain gap between the labeled and unlabeled target distributions and the inter-domain gap between source and unlabeled target distribution in SSDA. We suggest employing class-wise contrastive learning to reduce the inter-domain gap and instance-level contrastive alignment between the original(input image) and strongly augmented unlabeled target images to minimize the intra-domain discrepancy. We have empirically shown that both of these modules complement each other to achieve superior performance. Experiments on three well-known domain adaptation benchmark datasets, namely DomainNet, Office-Home, and Office31, demonstrate the effectiveness of our approach. CLDA achieves state-of-the-art results on all the above datasets.

Differentially Private n-gram Extraction

- Kunho Kim, Sivakanth Gopi, Janardhan Kulkarni, Sergey Yekhanin
- abstract@[open-review](#): We revisit the problem of $\$n$ -gram extraction in the differential privacy setting. In this problem, given a corpus of private text data, the goal is to release as many $\$n$ -grams as possible while preserving user level privacy. Extracting $\$n$ -grams is a fundamental subroutine in many NLP applications such as sentence completion, auto response generation for emails, etc. The problem also arises in other applications such as sequence mining, trajectory analysis, etc., and is a generalization of recently studied differentially private set union (DPSU) by Gopi et al. (2020). In this paper, we develop a new differentially private algorithm for this problem which, in our experiments, significantly outperforms the state-of-the-art. Our improvements stem from combining recent advances in DPSU, privacy accounting, and new heuristics for pruning in the tree-based approach initiated by Chen et al. (2012).

Capturing implicit hierarchical structure in 3D biomedical images with self-supervised hyperbolic representations

- Joy Hsu, Jeffrey Gu, Gong Wu, Wah Chiu, Serena Yeung
- abstract@[open-review](#): We consider the task of representation learning for unsupervised segmentation of 3D voxel-grid biomedical images. We show that models that capture implicit hierarchical relationships between subvolumes are better suited for this task. To that end, we consider encoder-decoder architectures with a hyperbolic latent space, to explicitly capture hierarchical relationships present in subvolumes of the data. We propose utilizing a 3D hyperbolic variational autoencoder with a novel gyroplane convolutional layer to map from the embedding space back to 3D images. To capture these relationships, we introduce an essential self-supervised loss---in addition to the standard VAE loss---which infers approximate hierarchies and encourages implicitly related subvolumes to be mapped closer in the embedding space. We present experiments on synthetic datasets along with a dataset from the medical domain to validate our hypothesis.

Noisy Recurrent Neural Networks

- Soon Hoe Lim, N. Benjamin Erichson, Liam Hodgkinson, Michael W. Mahoney
- abstract@[open-review](#): We provide a general framework for studying recurrent neural networks (RNNs) trained by injecting noise into hidden states. Specifically, we consider RNNs that can be viewed as discretizations of stochastic differential equations driven by input data. This framework allows us to study the implicit regularization effect of general noise injection schemes by deriving an approximate explicit regularizer in the small noise regime. We find that, under reasonable assumptions, this implicit regularization promotes flatter minima; it biases towards models with more stable dynamics; and, in classification tasks, it favors models with larger classification margin. Sufficient conditions for global stability are obtained, highlighting the phenomenon of stochastic stabilization, where noise injection can improve stability during training. Our theory is supported by empirical results which demonstrate that the RNNs have improved robustness with respect to various input perturbations.

Matrix encoding networks for neural combinatorial optimization

- Yeong-Dae Kwon, Jinho Choo, Iljoo Yoon, Minah Park, Duwon Park, Youngjune Gwon
- abstract@[open-review](#): Machine Learning (ML) can help solve combinatorial optimization (CO) problems better. A popular approach is to use a neural net to compute on the parameters of a given CO problem and extract useful information that guides the search for good solutions. Many CO problems of practical importance can be specified in a matrix form of parameters quantifying the relationship between two groups of items. There is currently no neural net model, however, that takes in such matrix-style relationship data as an input. Consequently, these types of CO problems have been out of reach

for ML engineers. In this paper, we introduce Matrix Encoding Network (MatNet) and show how conveniently it takes in and processes parameters of such complex CO problems. Using an end-to-end model based on MatNet, we solve asymmetric traveling salesman (ATSP) and flexible flow shop (FFSP) problems as the earliest neural approach. In particular, for a class of FFSP we have tested MatNet on, we demonstrate a far superior empirical performance to any methods (neural or not) known to date.

[When Is Unsupervised Disentanglement Possible?](#)

- Daniella Horan, Eitan Richardson, Yair Weiss
- abstract@[open-review](#): A common assumption in many domains is that high dimensional data are a smooth nonlinear function of a small number of independent factors. When is it possible to recover the factors from unlabeled data? In the context of deep models this problem is called “disentanglement” and was recently shown to be impossible without additional strong assumptions [17, 19]. In this paper, we show that the assumption of local isometry together with non-Gaussianity of the factors, is sufficient to provably recover disentangled representations from data. We leverage recent advances in deep generative models to construct manifolds of highly realistic images for which the ground truth latent representation is known, and test whether modern and classical methods succeed in recovering the latent factors. For many different manifolds, we find that a spectral method that explicitly optimizes local isometry and non-Gaussianity consistently finds the correct latent factors, while baseline deep autoencoders do not. We propose how to encourage deep autoencoders to find encodings that satisfy local isometry and show that this helps them discover disentangled representations. Overall, our results suggest that in some realistic settings, unsupervised disentanglement is provably possible, without any domain-specific assumptions.

[Continuous Latent Process Flows](#)

- Ruizhi Deng, Marcus A. Brubaker, Greg Mori, Andreas Lehrmann
- abstract@[open-review](#): Partial observations of continuous time-series dynamics at arbitrary time stamps exist in many disciplines. Fitting this type of data using statistical models with continuous dynamics is not only promising at an intuitive level but also has practical benefits, including the ability to generate continuous trajectories and to perform inference on previously unseen time stamps. Despite exciting progress in this area, the existing models still face challenges in terms of their representational power and the quality of their variational approximations. We tackle these challenges with continuous latent process flows (CLPF), a principled architecture decoding continuous latent processes into continuous observable processes using a time-dependent normalizing flow driven by a stochastic differential equation. To optimize our model using maximum likelihood, we propose a novel piecewise construction of a variational posterior process and derive the corresponding variational lower bound using trajectory re-weighting. Our ablation studies demonstrate the effectiveness of our contributions in various inference tasks on irregular time grids. Comparisons to state-of-the-art baselines show our model's favourable performance on both synthetic and real-world time-series data.

[Perturbation-based Regret Analysis of Predictive Control in Linear Time Varying Systems](#)

- Yiheng Lin, Yang Hu, Guanya Shi, Haoyuan Sun, Guannan Qu, Adam Wierman
- abstract@[open-review](#): We study predictive control in a setting where the dynamics are time-varying and linear, and the costs are time-varying and well-conditioned. At each time step, the controller receives the exact predictions of costs, dynamics, and disturbances for the future $\$k\$$ time steps. We show that when the prediction window $\$k\$$ is sufficiently large, predictive control is input-to-state stable and achieves a dynamic regret of $\$O(\lambda^k T)\$$, where $\lambda < 1$ is a positive constant. This is the first dynamic regret bound on the predictive control of linear time-varying systems. We also show a variation of predictive control obtains the first competitive bound for the control of linear time-varying systems: $\$1 + O(\lambda^k)\$$. Our results are derived using a novel proof framework based on a perturbation bound that characterizes how a small change to the system parameters impacts the optimal trajectory.

[Dataset Distillation with Infinitely Wide Convolutional Networks](#)

- Timothy Nguyen, Roman Novak, Lechao Xiao, Jaehoon Lee
- abstract@[open-review](#): The effectiveness of machine learning algorithms arises from being able to extract useful features from large amounts of data. As model and dataset sizes increase, dataset distillation methods that compress large datasets into significantly smaller yet highly performant ones will become valuable in terms of training efficiency and useful feature extraction. To that end, we apply a novel distributed kernel-based meta-learning framework to achieve state-of-the-art results for dataset distillation using infinitely wide convolutional neural networks. For instance, using only 10 datapoints (0.02% of original dataset), we obtain over 65% test accuracy on CIFAR-10 image classification task, a dramatic improvement over the previous best test accuracy of 40%. Our state-of-the-art results extend across many other settings for MNIST, Fashion-MNIST, CIFAR-10, CIFAR-100, and SVHN. Furthermore, we perform some preliminary analyses of our distilled datasets to shed light on how they differ from naturally occurring data.

[SPANN: Highly-efficient Billion-scale Approximate Nearest Neighborhood Search](#)

- Qi Chen, Bing Zhao, Haidong Wang, Mingqin Li, Chuanjie Liu, Zengzhong Li, Mao Yang, Jingdong Wang
- abstract@[open-review](#): The in-memory algorithms for approximate nearest neighbor search (ANNS) have achieved great success for fast high-recall search, but are extremely expensive when handling very large scale database. Thus, there is an increasing request for the hybrid ANNS solutions with small memory and inexpensive solid-state drive (SSD). In this paper, we present a simple but efficient memory-disk hybrid indexing and search system, named SPANN, that follows the inverted index methodology. It stores the centroid points of the posting lists in the memory and the large posting lists in the disk. We guarantee both disk-access efficiency (low latency) and high recall by effectively reducing the disk-access number and retrieving high-quality posting lists. In the index-building stage, we adopt a hierarchical balanced clustering algorithm to balance the length of posting lists and augment the posting list by adding the points in the closure of the corresponding clusters. In the search stage, we use a query-aware scheme to dynamically prune the access of unnecessary posting lists. Experiment results demonstrate that SPANN is 2X faster than the state-of-the-art ANNS solution DiskANN to reach the same recall quality 90% with same memory cost in three billion-scale datasets. It can reach 90% recall@1 and recall@10 in just around one millisecond with only about 10% of original memory cost. Code is available at: <https://github.com/microsoft/SPTAG>.

[Distilling Object Detectors with Feature Richness](#)

- Du Zhixing, Rui Zhang, Ming Chang, xishan zhang, Shaoli Liu, Tianshi Chen, Yunji Chen
- abstract@[open-review](#): In recent years, large-scale deep models have achieved great success, but the huge computational complexity and massive storage requirements make it a great challenge to deploy them in resource-limited devices. As a model compression and acceleration method, knowledge distillation effectively improves the performance of small models by transferring the dark knowledge from the teacher detector. However, most of the existing distillation-based detection methods mainly imitating features near bounding boxes, which suffer from two limitations. First, they ignore the beneficial features outside the bounding boxes. Second, these methods imitate some features which are mistakenly regarded as the background by the teacher detector. To address the above issues, we propose a novel Feature-Richness Score (FRS) method to choose important features that improve generalized detectability during distilling. The proposed method effectively retrieves the important features outside the bounding boxes and removes the detrimental features within the bounding boxes. Extensive experiments show that our methods achieve excellent performance on both anchor-based and anchor-free detectors. For example, RetinaNet with ResNet-50 achieves 39.7% in mAP on the COCO2017 dataset, which even surpasses the ResNet-101 based teacher detector 38.9% by 0.8%. Our implementation is available at <https://github.com/duzhixing/FRS>.

[Analysis of one-hidden-layer neural networks via the resolvent method](#)

- Vanessa Piccolo, Dominik Schröder
- abstract@[open-review](#): In this work, we investigate the asymptotic spectral density of the random feature matrix $M = Y Y^*$ with $Y = f(WX)$ generated by a single-hidden-layer neural network, where W and X are random rectangular matrices with i.i.d. centred entries and f is a non-linear smooth function which is applied entry-wise. We prove that the Stieltjes transform of the limiting spectral distribution approximately satisfies a quartic self-consistent equation, which is exactly the equation obtained by [Pennington, Worah 2017] and [Benigni, Péché 2019] with the moment method. We extend the previous results to the case of additive bias $Y = f(WX + B)$ with B being an independent rank-one Gaussian random matrix, closer modelling the neural network infrastructures encountered in practice. Our key finding is that in the case of additive bias it is impossible to choose an activation function preserving the layer-to-layer singular value distribution, in sharp contrast to the bias-free case where a simple integral constraint is sufficient to achieve isospectrality. To obtain the asymptotics for the empirical spectral density we follow the resolvent method from random matrix theory via the cumulant expansion. We find that this approach is more robust and less combinatorial than the moment method and expect that it will apply also for models where the combinatorics of the former become intractable. The resolvent method has been widely employed, but compared to previous works, it is applied here to non-linear random matrices.

[Grounding Spatio-Temporal Language with Transformers](#)

- Tristan Karch, Laetitia Teodorescu, Katja Hofmann, Clément Moulin-Frier, Pierre-Yves Oudeyer
- abstract@[open-review](#): Language is an interface to the outside world. In order for embodied agents to use it, language must be grounded in other, sensorimotor modalities. While there is an extended literature studying how machines can learn grounded language, the topic of how to learn spatio-temporal linguistic concepts is still largely uncharted. To make progress in this direction, we here introduce a novel spatio-temporal language grounding task where the goal is to learn the meaning of spatio-temporal descriptions of behavioral traces of an embodied agent. This is achieved by training a truth function that predicts if a description matches a given history of observations. The descriptions involve time-extended predicates in past and present tense as well as spatio-temporal references to objects in the scene. To study the role of architectural biases in this task, we train several models including multimodal Transformer architectures; the latter implement different attention computations between words and objects across space and time. We test models on two classes of generalization: 1) generalization to new sentences, 2) generalization to grammar primitives. We observe that maintaining object identity in the attention computation of our Transformers is instrumental to achieving good performance on generalization overall, and that summarizing object traces in a single token has little influence on performance. We then discuss how this opens new perspectives for language-guided autonomous embodied agents.

[Learning where to learn: Gradient sparsity in meta and continual learning](#)

- Johannes von Oswald, Dominic Zhao, Seijin Kobayashi, Simon Schug, Massimo Caccia, Nicolas Zucchet, João Sacramento
- abstract@[open-review](#): Finding neural network weights that generalize well from small datasets is difficult. A promising approach is to learn a weight initialization such that a small number of weight changes results in low generalization error. We show that this form of meta-learning can be improved by letting the learning algorithm decide which weights to change, i.e., by learning where to learn. We find that patterned sparsity emerges from this process, with the pattern of sparsity varying on a problem-by-problem basis. This selective sparsity results in better generalization and less interference in a range of few-shot and continual learning problems. Moreover, we find that sparse learning also emerges in a more expressive model where learning rates are meta-learned. Our results shed light on an ongoing debate on whether meta-learning can discover adaptable features and suggest that learning by sparse gradient descent is a powerful inductive bias for meta-learning systems.

[Domain Invariant Representation Learning with Domain Density Transformations](#)

- A. Tuan Nguyen, Toan Tran, Yarin Gal, Atilim Gunes Baydin
- abstract@[open-review](#): Domain generalization refers to the problem where we aim to train a model on data from a set of source domains so that the model can generalize to unseen target domains. Naively training a model on the aggregate set of data (pooled from all source domains) has been shown to perform suboptimally, since the information learned by that model might be domain-specific and generalize imperfectly to target domains. To tackle this problem, a predominant domain generalization approach is to learn some domain-invariant information for the prediction task, aiming at a good generalization across domains. In this paper, we propose a theoretically grounded method to learn a domain-invariant representation by enforcing the representation network to be invariant under all transformation functions among domains. We next introduce the use of generative adversarial networks to learn such domain transformations in a possible implementation of our method in practice. We demonstrate the effectiveness of our method on several widely used datasets for the domain generalization problem, on all of which we achieve competitive results with state-of-the-art models.

[PlayVirtual: Augmenting Cycle-Consistent Virtual Trajectories for Reinforcement Learning](#)

- Tao Yu, Cuiling Lan, Wenjun Zeng, Mingxiao Feng, Zhizheng Zhang, Zhibo Chen
- abstract@[open-review](#): Learning good feature representations is important for deep reinforcement learning (RL). However, with limited experience, RL often suffers from data inefficiency for training. For un-experienced or less-experienced trajectories (i.e., state-action sequences), the lack of data limits the use of them for better feature learning. In this work, we propose a novel method, dubbed PlayVirtual, which augments cycle-consistent virtual trajectories to enhance the data efficiency for RL feature representation learning. Specifically, PlayVirtual predicts future states in a latent space based on the current state and action by a dynamics model and then predicts the previous states by a backward dynamics model, which forms a trajectory cycle. Based on this, we augment the actions to generate a large amount of virtual state-action trajectories. Being free of groundtruth state supervision, we enforce a trajectory to meet the cycle consistency constraint, which can significantly enhance the data efficiency. We validate the effectiveness of our designs on the Atari and DeepMind Control Suite benchmarks. Our method achieves the state-of-the-art performance on both benchmarks. Our code is available at <https://github.com/microsoft/Playvirtual>.

[Efficient Equivariant Network](#)

- Lingshen He, Yuxuan Chen, zhengyang shen, Yiming Dong, Yisen Wang, Zhouchen Lin
- abstract@[open-review](#): Convolutional neural networks (CNNs) have dominated the field of Computer Vision and achieved great success due to their built-in translation equivariance. Group equivariant CNNs (G-CNNs) that incorporate more equivariance can significantly improve the performance of conventional CNNs. However, G-CNNs are faced with two major challenges: \{spatial-agnostic problem\} and \{expensive computational cost\}. In this work, we propose a general framework of previous equivariant models, which includes G-CNNs and equivariant self-attention layers as special cases. Under this framework, we explicitly decompose the feature aggregation operation into a kernel generator and an encoder, and decouple the spatial and extra geometric dimensions in the computation. Therefore, our filters are essentially dynamic rather than being spatial-agnostic. We further show that our \{E\}quivariant model is parameter \{E\}fficient and computation \{E\}fficient by complexity analysis, and also data \{E\}fficient by experiments, so we call our model E^4 -Net. Extensive experiments verify that our model can significantly improve previous works with smaller model size. Especially, under the setting of training on $1/5$ data of CIFAR10, our model improves G-CNNs by $5\%+$ accuracy, while using only 56% parameters and 68% FLOPs.

[Unifying Gradient Estimators for Meta-Reinforcement Learning via Off-Policy Evaluation](#)

- Yunhao Tang, Tadashi Kozuno, Mark Rowland, Remi Munos, Michal Valko
- abstract@[open-review](#): Model-agnostic meta-reinforcement learning requires estimating the Hessian matrix of value functions. This is challenging from an implementation perspective, as repeatedly differentiating policy gradient estimates may lead to biased Hessian estimates. In this work, we provide a unifying framework for estimating higher-order derivatives of value functions, based on off-policy evaluation. Our framework interprets a number of prior approaches as special cases and elucidates the bias and variance trade-off of Hessian estimates. This framework also opens the door to a new family of estimates, which can be easily implemented with auto-differentiation libraries, and lead to performance gains in practice.

[Even your Teacher Needs Guidance: Ground-Truth Targets Dampen Regularization Imposed by Self-Distillation](#)

- Kenneth Borup, Lars N Andersen
- abstract@[open-review](#): Knowledge distillation is classically a procedure where a neural network is trained on the output of another network along with the original targets in order to transfer knowledge between the architectures. The special case of self-distillation, where the network architectures are identical, has been observed to improve generalization accuracy. In this paper, we consider an iterative variant of self-distillation in a kernel regression setting, in which successive steps incorporate both model outputs and the ground-truth targets. This allows us to provide the first theoretical results on the importance of using the weighted ground-truth targets in self-distillation. Our focus is on fitting nonlinear functions to training data with a weighted mean square error objective function suitable for distillation, subject to $\|\cdot\|_2$ regularization of the model parameters. We show that any such function obtained with self-distillation can be calculated directly as a function of the initial fit, and that infinite distillation steps yields the same optimization problem as the original with amplified regularization. Furthermore, we provide a closed form solution for the optimal choice of weighting parameter at each step, and show how to efficiently estimate this weighting parameter for deep learning and significantly reduce the computational requirements compared to a grid search.

[Compressing Neural Networks: Towards Determining the Optimal Layer-wise Decomposition](#)

- Lucas Liebenwein, Alaa Maalouf, Dan Feldman, Daniela Rus
- abstract@[open-review](#): We present a novel global compression framework for deep neural networks that automatically analyzes each layer to identify the optimal per-layer compression ratio, while simultaneously achieving the desired overall compression. Our algorithm hinges on the idea of compressing each convolutional (or fully-connected) layer by slicing its channels into multiple groups and decomposing each group via low-rank decomposition. At the core of our algorithm is the derivation of layer-wise error bounds from the Eckart–Young–Mirsky theorem. We then leverage these bounds to frame the compression problem as an optimization problem where we wish to minimize the maximum compression error across layers and propose an efficient algorithm towards a solution. Our experiments indicate that our method outperforms existing low-rank compression approaches across a wide range of networks and data sets. We believe that our results open up new avenues for future research into the global performance-size trade-offs of modern neural networks.

[Equilibrium and non-Equilibrium regimes in the learning of Restricted Boltzmann Machines](#)

- Aurélien Decelle, Cyril Furtlehner, Beatriz Seoane
- abstract@[open-review](#): Training Restricted Boltzmann Machines (RBMs) has been challenging for a long time due to the difficulty of computing precisely the log-likelihood gradient. Over the past decades, many works have proposed more or less successful recipes but without studying systematically the crucial quantity of the problem: the mixing time i.e. the number of MCMC iterations needed to sample completely new configurations from a model. In this work, we show that this mixing time plays a crucial role in the behavior and stability of the trained model, and that RBMs operate in two well-defined distinct regimes, namely equilibrium and out-of-equilibrium, depending on the interplay between this mixing time of the model and the number of MCMC steps, k , used to approximate the gradient. We further show empirically that this mixing time increases along the learning, which often implies a transition from one regime to another as soon as k becomes smaller than this time. In particular, we show that using the popular k (persistent) contrastive divergence approaches, with k small, the dynamics of the fitted model are extremely slow and often dominated by strong out-of-equilibrium effects. On the contrary, RBMs trained in equilibrium display much faster dynamics, and a smooth convergence to dataset-like configurations during the sampling. Finally, we discuss how to exploit in practice both regimes depending on the task one aims to fulfill: (i) short k s can be used to generate convincing samples in short learning times, (ii) large k s (or increasingly large) must be used to learn the correct equilibrium distribution of the RBM. Finally, the existence of these two operational regimes seems to be a general property of energy based models trained via likelihood maximization.

[Imitation with Neural Density Models](#)

- Kuno Kim, Akshat Jindal, Yang Song, Jiaming Song, Yanan Sui, Stefano Ermon
- abstract@[open-review](#): We propose a new framework for Imitation Learning (IL) via density estimation of the expert's occupancy measure followed by Maximum Occupancy Entropy Reinforcement Learning (RL) using the density as a reward. Our approach maximizes a non-adversarial model-free RL objective that provably lower bounds reverse Kullback–Leibler divergence between occupancy measures of the expert and imitator. We present a practical IL algorithm, Neural Density Imitation (NDI), which obtains state-of-the-art demonstration efficiency on benchmark control tasks.

[Accurate Point Cloud Registration with Robust Optimal Transport](#)

- Zhengyang Shen, Jean Feydy, Peirong Liu, Ariel H Curiale, Ruben San Jose Estepar, Raul San Jose Estepar, Marc Niethammer
- abstract@[open-review](#): This work investigates the use of robust optimal transport (OT) for shape matching. Specifically, we show that recent OT solvers improve both optimization-based and deep learning methods for point cloud registration, boosting accuracy at an affordable computational cost. This manuscript starts with a practical overview of modern OT theory. We then provide solutions to the main difficulties in using this framework for shape matching. Finally, we showcase the performance of transport-enhanced registration models on a wide range of challenging tasks: rigid registration for partial shapes; scene flow estimation on the Kitti dataset; and nonparametric registration of lung vascular trees between inspiration and expiration. Our OT-based methods achieve state-of-the-art results on Kitti and for the challenging lung registration task, both in terms of accuracy and scalability. We also release PVT1010, a new public dataset of 1,010 pairs of lung vascular trees with densely sampled points. This dataset provides a challenging use case for point cloud registration algorithms with highly complex shapes and deformations. Our work demonstrates that robust OT enables fast pre-alignment and fine-tuning for a wide range of registration models, thereby providing a new key method for the computer vision toolbox. Our code and dataset are available online at: <https://github.com/uncbiag/robot>.

[Simple steps are all you need: Frank-Wolfe and generalized self-concordant functions](#)

- Alejandro Carderera, Mathieu Besançon, Sebastian Pokutta
- abstract@[open-review](#): Generalized self-concordance is a key property present in the objective function of many important learning problems. We establish the convergence rate of a simple Frank-Wolfe variant that uses the open-loop step size strategy $\gamma_t = 2/(t+2)$, obtaining a $(1/t)$ convergence rate for this class of functions in terms of primal gap and Frank-Wolfe gap, where t is the iteration count. This avoids the use of second-order information or the need to estimate local smoothness parameters of previous work. We also show improved convergence rates for various common cases, e.g., when the feasible region under consideration is uniformly convex or polyhedral.

[Automatic Data Augmentation for Generalization in Reinforcement Learning](#)

- Roberta Raileanu, Maxwell Goldstein, Denis Yarats, Ilya Kostrikov, Rob Fergus
- abstract@[open-review](#): Deep reinforcement learning (RL) agents often fail to generalize beyond their training environments. To alleviate this problem, recent work has proposed the use of data augmentation. However, different tasks tend to benefit from different types of augmentations and selecting the right one typically requires expert knowledge. In this paper, we introduce three approaches for automatically finding an effective augmentation for any RL task. These are combined with two novel regularization terms for the policy and value function, required to make the use of data augmentation theoretically sound for actor-critic algorithms. Our method achieves a new state-of-the-art on the Procgen benchmark and outperforms popular RL algorithms on DeepMind Control tasks with distractors. In addition, our agent learns policies and representations which are more robust to changes in the environment that are irrelevant for solving the task, such as the background.

[Blending Anti-Aliasing into Vision Transformer](#)

- Shengju Qian, Hao Shao, Yi Zhu, Mu Li, Jiaya Jia
- abstract@[open-review](#): The transformer architectures, based on self-attention mechanism and convolution-free design, recently found superior performance and booming applications in computer vision. However, the discontinuous patch-wise tokenization process implicitly introduces jagged artifacts into attention maps, arising the traditional problem of aliasing for vision transformers. Aliasing effect occurs when discrete patterns are used to produce high frequency or continuous information, resulting in the indistinguishable distortions. Recent researches have found that modern convolution networks still suffer from this phenomenon. In this work, we analyze the uncharted problem of aliasing in vision transformer and explore to incorporate anti-aliasing properties. Specifically, we propose a plug-and-play Aliasing-Reduction Module (ARM) to alleviate the aforementioned issue. We investigate the effectiveness and generalization of the proposed method across multiple tasks and various vision transformer families. This lightweight design consistently attains a clear boost over several famous structures. Furthermore, our module also improves data efficiency and robustness of vision transformers.

[A Trainable Spectral-Spatial Sparse Coding Model for Hyperspectral Image Restoration](#)

- Theo Bodro, Alexandre Zouaoui, Jocelyn Chanussot, Julien Mairal
- abstract@[open-review](#): Hyperspectral imaging offers new perspectives for diverse applications, ranging from the monitoring of the environment using airborne or satellite remote sensing, precision farming, food safety, planetary exploration, or astrophysics. Unfortunately, the spectral diversity of information comes at the expense of various sources of degradation, and the lack of accurate ground-truth "clean" hyperspectral signals acquired on the spot makes restoration tasks challenging. In particular, training deep neural networks for restoration is difficult, in contrast to traditional RGB imaging problems where deep models tend to shine. In this paper, we advocate instead for a hybrid approach based on sparse coding principles that retain the interpretability of classical techniques encoding domain knowledge with handcrafted image priors, while allowing to train model parameters end-to-end without massive amounts of data. We show on various denoising benchmarks that our method is computationally efficient and significantly outperforms the state of the art.

[Posterior Collapse and Latent Variable Non-identifiability](#)

- Yixin Wang, David Blei, John P. Cunningham
- abstract@[open-review](#): Variational autoencoders model high-dimensional data by positing low-dimensional latent variables that are mapped through a flexible distribution parametrized by a neural network. Unfortunately, variational autoencoders often suffer from posterior collapse: the posterior of the latent variables is equal to its prior, rendering the variational autoencoder useless as a means to produce meaningful representations. Existing approaches to posterior collapse often attribute it to the use of neural networks or optimization issues due to variational approximation. In this paper, we consider posterior collapse as a problem of latent variable non-identifiability. We prove that the posterior collapses if and only if the latent variables are non-identifiable in the generative model. This fact implies that posterior collapse is not a phenomenon specific to the use of flexible distributions or approximate inference. Rather, it can occur in classical probabilistic models even with exact inference, which we also demonstrate. Based on these results, we propose a class of latent-identifiable variational autoencoders, deep generative models which enforce identifiability without sacrificing flexibility. This model class resolves the problem of latent variable non-identifiability by leveraging bijective Brenier maps and parameterizing them with input convex neural networks, without special variational inference objectives or optimization tricks. Across synthetic and real datasets, latent-identifiable variational autoencoders outperform existing methods in mitigating posterior collapse and providing meaningful representations of the data.

[The Benefits of Implicit Regularization from SGD in Least Squares Problems](#)

- Difan Zou, Jingfeng Wu, Vladimir Braverman, Quanquan Gu, Dean P. Foster, Sham Kakade
- abstract@[open-review](#): Stochastic gradient descent (SGD) exhibits strong algorithmic regularization effects in practice, which has been hypothesized to play an important role in the generalization of modern machine learning approaches. In this work, we seek to understand these issues in the simpler setting of linear regression (including both underparameterized and overparameterized regimes), where our goal is to make sharp instance-based comparisons of the implicit regularization afforded by (unregularized) average SGD with the explicit regularization of ridge regression. For a broad class of least squares problem instances (that are natural in high-dimensional settings), we show: (1) for every problem instance and for every ridge parameter, (unregularized) SGD, when provided with logarithmically more samples than that provided to the ridge algorithm, generalizes no worse than the ridge solution (provided SGD uses a tuned constant stepsize); (2) conversely, there exist instances (in this wide problem class) where optimally-tuned ridge regression requires quadratically more samples than SGD in order to have the same generalization performance. Taken together, our results show that, up to the logarithmic factors, the generalization performance of SGD is always no worse than that of ridge regression in a wide range of overparameterized problems, and, in fact, could be much better for some problem instances. More generally, our results show how algorithmic regularization has important consequences even in simpler (overparameterized) convex settings.

[Generalization of Model-Agnostic Meta-Learning Algorithms: Recurring and Unseen Tasks](#)

- Alireza Fallah, Aryan Mokhtari, Asuman Ozdaglar
- abstract@[open-review](#): In this paper, we study the generalization properties of Model-Agnostic Meta-Learning (MAML) algorithms for supervised learning problems. We focus on the setting in which we train the MAML model over m tasks, each with n data points, and characterize its generalization error from two points of view: First, we assume the new task at test time is one of the training tasks, and we show that, for strongly convex objective functions, the expected excess population loss is bounded by $O(1/mn)$. Second, we consider the MAML algorithm's generalization to an unseen task and show that the resulting generalization error depends on the total variation distance between the underlying distributions of the new task and the tasks observed during the training process. Our proof techniques rely on the connections between algorithmic stability and generalization bounds of algorithms. In particular, we propose a new definition of stability for meta-learning algorithms, which allows us to capture the role of both the number of tasks m and number of samples per task n on the generalization error of MAML.

[Factored Policy Gradients: Leveraging Structure for Efficient Learning in MDPs](#)

- Thomas Spooner, Nelson Vadovi, Sumitra Ganesh
- abstract@[open-review](#): Policy gradient methods can solve complex tasks but often fail when the dimensionality of the action-space or objective multiplicity grow very large. This occurs, in part, because the variance on score-based gradient estimators scales quadratically. In this paper, we address this problem through a factor baseline which exploits independence structure encoded in a novel action-target influence network. Factored policy gradients

(FPGs), which follow, provide a common framework for analysing key state-of-the-art algorithms, are shown to generalise traditional policy gradients, and yield a principled way of incorporating prior knowledge of a problem domain's generative processes. We provide an analysis of the proposed estimator and identify the conditions under which variance is reduced. The algorithmic aspects of FPGs are discussed, including optimal policy factorisation, as characterised by minimum biclique coverings, and the implications for the bias variance trade-off of incorrectly specifying the network. Finally, we demonstrate the performance advantages of our algorithm on large-scale bandit and traffic intersection problems, providing a novel contribution to the latter in the form of a spatial approximation.

[MarioNette: Self-Supervised Sprite Learning](#)

- Dmitriy Smirnov, MICHAEL GHARBI, Matthew Fisher, Vitor Guizilini, Alexei Efros, Justin M. Solomon
- abstract@[open-review](#): Artists and video game designers often construct 2D animations using libraries of sprites---textured patches of objects and characters. We propose a deep learning approach that decomposes sprite-based video animations into a disentangled representation of recurring graphic elements in a self-supervised manner. By jointly learning a dictionary of possibly transparent patches and training a network that places them onto a canvas, we deconstruct sprite-based content into a sparse, consistent, and explicit representation that can be easily used in downstream tasks, like editing or analysis. Our framework offers a promising approach for discovering recurring visual patterns in image collections without supervision.

[RLlib Flow: Distributed Reinforcement Learning is a Dataflow Problem](#)

- Eric Liang, Zhanghao Wu, Michael Luo, Sven Mika, Joseph E. Gonzalez, Ion Stoica
- abstract@[open-review](#): Researchers and practitioners in the field of reinforcement learning (RL) frequently leverage parallel computation, which has led to a plethora of new algorithms and systems in the last few years. In this paper, we re-examine the challenges posed by distributed RL and try to view it through the lens of an old idea: distributed dataflow. We show that viewing RL as a dataflow problem leads to highly composable and performant implementations. We propose RLlib Flow, a hybrid actor-dataflow programming model for distributed RL, and validate its practicality by porting the full suite of algorithms in RLlib, a widely adopted distributed RL library. Concretely, RLlib Flow provides 2-9\$times\$ code savings in real production code and enables the composition of multi-agent algorithms not possible by end users before. The open-source code is available as part of RLlib at <https://github.com/ray-project/ray/tree/master/rllib>.

[Improve Agents without Retraining: Parallel Tree Search with Off-Policy Correction](#)

- Gal Dalal, Assaf Hallak, Steven Dalton, iuri frosio, Shie Mannor, Gal Chechik
- abstract@[open-review](#): Tree Search (TS) is crucial to some of the most influential successes in reinforcement learning. Here, we tackle two major challenges with TS that limit its usability: \textit{distribution shift} and \textit{scalability}. We first discover and analyze a counter-intuitive phenomenon: action selection through TS and a pre-trained value function often leads to lower performance compared to the original pre-trained agent, even when having access to the exact state and reward in future steps. We show this is due to a distribution shift to areas where value estimates are highly inaccurate and analyze this effect using Extreme Value theory. To overcome this problem, we introduce a novel off-policy correction term that accounts for the mismatch between the pre-trained value and its corresponding TS policy by penalizing under-sampled trajectories. We prove that our correction eliminates the above mismatch and bound the probability of sub-optimal action selection. Our correction significantly improves pre-trained Rainbow agents without any further training, often more than doubling their scores on Atari games. Next, we address the scalability issue given by the computational complexity of exhaustive TS that scales exponentially with the tree depth. We introduce Batch-BFS: a GPU breadth-first search that advances all nodes in each depth of the tree simultaneously. Batch-BFS reduces runtime by two orders of magnitude and, beyond inference, enables also training with TS of depths that were not feasible before. We train DQN agents from scratch using TS and show improvement in several Atari games compared to both the original DQN and the more advanced Rainbow. We will share the code upon publication.

[Redesigning the Transformer Architecture with Insights from Multi-particle Dynamical Systems](#)

- Subhabrata Dutta, Tanya Gautam, Soumen Chakrabarti, Tanmoy Chakraborty
- abstract@[open-review](#): The Transformer and its variants have been proven to be efficient sequence learners in many different domains. Despite their staggering success, a critical issue has been the enormous number of parameters that must be trained (ranging from 10^7 to 10^{11}) along with the quadratic complexity of dot-product attention. In this work, we investigate the problem of approximating the two central components of the Transformer -- multi-head self-attention and point-wise feed-forward transformation, with reduced parameter space and computational complexity. We build upon recent developments in analyzing deep neural networks as numerical solvers of ordinary differential equations. Taking advantage of an analogy between Transformer stages and the evolution of a dynamical system of multiple interacting particles, we formulate a temporal evolution scheme, \name, to bypass costly dot-product attention over multiple stacked layers. We perform exhaustive experiments with \name\ on well-known encoder-decoder as well as encoder-only tasks. We observe that the degree of approximation (or inversely, the degree of parameter reduction) has different effects on the performance, depending on the task. While in the encoder-decoder regime, \name\ delivers performances comparable to the original Transformer, in encoder-only tasks it consistently outperforms Transformer along with several subsequent variants.

[Exploring Architectural Ingredients of Adversarially Robust Deep Neural Networks](#)

- Hanxun Huang, Yisen Wang, Sarah Erfani, Quanquan Gu, James Bailey, Xingjun Ma
- abstract@[open-review](#): Deep neural networks (DNNs) are known to be vulnerable to adversarial attacks. A range of defense methods have been proposed to train adversarially robust DNNs, among which adversarial training has demonstrated promising results. However, despite preliminary understandings developed for adversarial training, it is still not clear, from the architectural perspective, what configurations can lead to more robust DNNs. In this paper, we address this gap via a comprehensive investigation on the impact of network width and depth on the robustness of adversarially trained DNNs. Specifically, we make the following key observations: 1) more parameters (higher model capacity) does not necessarily help adversarial robustness; 2) reducing capacity at the last stage (the last group of blocks) of the network can actually improve adversarial robustness; and 3) under the same parameter budget, there exists an optimal architectural configuration for adversarial robustness. We also provide a theoretical analysis explaining why such network configuration can help robustness. These architectural insights can help design adversarially robust DNNs.

[Center Smoothing: Certified Robustness for Networks with Structured Outputs](#)

- Aounon Kumar, Tom Goldstein
- abstract@[open-review](#): The study of provable adversarial robustness has mostly been limited to classification tasks and models with one-dimensional real-valued outputs. We extend the scope of certifiable robustness to problems with more general and structured outputs like sets, images, language, etc. We model the output space as a metric space under a distance/similarity function, such as intersection-over-union, perceptual similarity, total variation distance, etc. Such models are used in many machine learning problems like image segmentation, object detection, generative models, image/audio-to-text systems, etc. Based on a robustness technique called randomized smoothing, our center smoothing procedure can produce models with the guarantee that the change in the output, as measured by the distance metric, remains small for any norm-bounded adversarial perturbation of the input. We apply our method to create certifiably robust models with disparate output spaces -- from sets to images -- and show that it yields meaningful certificates without significantly degrading the performance of the base model.

[Breaking the Linear Iteration Cost Barrier for Some Well-known Conditional Gradient Methods Using MaxIP Data-structures](#)

- Zhaozhuo Xu, Zhao Song, Anshumali Shrivastava
- abstract@[open-review](#): Conditional gradient methods (CGM) are widely used in modern machine learning. CGM's overall running time usually consists of two parts: the number of iterations and the cost of each iteration. Most efforts focus on reducing the number of iterations as a means to reduce the overall running time. In this work, we focus on improving the per iteration cost of CGM. The bottleneck step in most CGM is maximum inner product search (MaxIP), which requires a linear scan over the parameters. In practice, approximate MaxIP data-structures are found to be helpful heuristics. However, theoretically, nothing is known about the combination of approximate MaxIP data-structures and CGM. In this work, we answer this question positively by providing a formal framework to combine the locality sensitive hashing type approximate MaxIP data-structures with CGM algorithms. As a result, we show the first algorithm, where the cost per iteration is sublinear in the number of parameters, for many fundamental optimization algorithms, e.g., Frank-Wolfe, Herding algorithm, and policy gradient.

[Neural Regression, Representational Similarity, Model Zoology & Neural Taskonomy at Scale in Rodent Visual Cortex](#)

- Colin Conwell, David Mayo, Andrei Barbu, Michael Buice, George Alvarez, Boris Katz
- abstract@[open-review](#): How well do deep neural networks fare as models of mouse visual cortex? A majority of research to date suggests results far more mixed than those produced in the modeling of primate visual cortex. Here, we perform a large-scale benchmarking of dozens of deep neural network models in mouse visual cortex with both representational similarity analysis and neural regression. Using the Allen Brain Observatory's 2-photon calcium-imaging dataset of activity in over 6,000 reliable rodent visual cortical neurons recorded in response to natural scenes, we replicate previous findings and resolve previous discrepancies, ultimately demonstrating that modern neural networks can in fact be used to explain activity in the mouse visual cortex to a more reasonable degree than previously suggested. Using our benchmark as an atlas, we offer preliminary answers to overarching questions about levels of analysis (e.g. do models that better predict the representations of individual neurons also predict representational similarity across neural populations?); questions about the properties of models that best predict the visual system overall (e.g. is convolution or category-supervision necessary to better predict neural activity?); and questions about the mapping between biological and artificial representations (e.g. does the information processing hierarchy in deep nets match the anatomical hierarchy of mouse visual cortex?). Along the way, we catalogue a number of models (including vision transformers, MLP-Mixers, normalization free networks, Taskonomy encoders and self-supervised models) outside the traditional circuit of convolutional object recognition. Taken together, our results provide a reference point for future ventures in the deep neural network modeling of mouse visual cortex, hinting at novel combinations of mapping method, architecture, and task to more fully characterize the computational motifs of visual representation in a species so central to neuroscience, but with a perceptual physiology and ecology markedly different from the ones we study in primates.

[A Topological Perspective on Causal Inference](#)

- Duligur Ibeling, Thomas Icard
- abstract@[open-review](#): This paper presents a topological learning-theoretic perspective on causal inference by introducing a series of topologies defined on general spaces of structural causal models (SCMs). As an illustration of the framework we prove a topological causal hierarchy theorem, showing that substantive assumption-free causal inference is possible only in a meager set of SCMs. Thanks to a known correspondence between open sets in the weak topology and statistically verifiable hypotheses, our results show that inductive assumptions sufficient to license valid causal inferences are statistically unverifiable in principle. Similar to no-free-lunch theorems for statistical inference, the present results clarify the inevitability of substantial assumptions for causal inference. An additional benefit of our topological approach is that it easily accommodates SCMs with infinitely many variables. We finally suggest that our framework may be helpful for the positive project of exploring and assessing alternative causal-inductive assumptions.

[Parameter Inference with Bifurcation Diagrams](#)

- Gregory Szep, Neil Dalchau, Attila Csikász-Nagy
- abstract@[open-review](#): Estimation of parameters in differential equation models can be achieved by applying learning algorithms to quantitative time-series data. However, sometimes it is only possible to measure qualitative changes of a system in response to a controlled condition. In dynamical systems theory, such change points are known as bifurcations and lie on a function of the controlled condition called the bifurcation diagram. In this work, we propose a gradient-based approach for inferring the parameters of differential equations that produce a user-specified bifurcation diagram. The cost function contains an error term that is minimal when the model bifurcations match the specified targets and a bifurcation measure which has gradients that push optimisers towards bifurcating parameter regimes. The gradients can be computed without the need to differentiate through the operations of the solver that was used to compute the diagram. We demonstrate parameter inference with minimal models which explore the space of saddle-node and pitchfork diagrams and the genetic toggle switch from synthetic biology. Furthermore, the cost landscape allows us to organise models in terms of topological and geometric equivalence.

[Scalable Thompson Sampling using Sparse Gaussian Process Models](#)

- Sattar Vakili, Henry Moss, Artem Artemev, Vincent Dutordoir, Victor Picheny
- abstract@[open-review](#): Thompson Sampling (TS) from Gaussian Process (GP) models is a powerful tool for the optimization of black-box functions. Although TS enjoys strong theoretical guarantees and convincing empirical performance, it incurs a large computational overhead that scales polynomially with the optimization budget. Recently, scalable TS methods based on sparse GP models have been proposed to increase the scope of TS, enabling its application to problems that are sufficiently multi-modal, noisy or combinatorial to require more than a few hundred evaluations to be solved. However, the approximation error introduced by sparse GPs invalidates all existing regret bounds. In this work, we perform a theoretical and empirical analysis of scalable TS. We provide theoretical guarantees and show that the drastic reduction in computational complexity of scalable TS can be enjoyed without loss in the regret performance over the standard TS. These conceptual claims are validated for practical implementations of scalable TS on synthetic benchmarks and as part of a real-world high-throughput molecular design task.

[Robust Counterfactual Explanations on Graph Neural Networks](#)

- Mohit Bajaj, Lingyang Chu, Zi Yu Xue, Jian Pei, Lanjun Wang, Peter Cho-Ho Lam, Yong Zhang
- abstract@[open-review](#): Massive deployment of Graph Neural Networks (GNNs) in high-stake applications generates a strong demand for explanations that are robust to noise and align well with human intuition. Most existing methods generate explanations by identifying a subgraph of an input graph that has a strong correlation with the prediction. These explanations are not robust to noise because independently optimizing the correlation for a single input can easily overfit noise. Moreover, they are not counterfactual because removing an identified subgraph from an input graph does not necessarily change the prediction result. In this paper, we propose a novel method to generate robust counterfactual explanations on GNNs by explicitly modelling the common decision logic of GNNs on similar input graphs. Our explanations are naturally robust to noise because they are produced from the common decision boundaries of a GNN that govern the predictions of many similar input graphs. The explanations are also counterfactual because removing the set of edges identified by an explanation from the input graph changes the prediction significantly. Exhaustive experiments on many public datasets demonstrate the superior performance of our method.

[Similarity and Matching of Neural Network Representations](#)

- Adrián Csiszárík, Péter Kőrösi-Szabó, Ákos Matsangosz, Gergely Papp, Dániel Varga
- abstract@[open-review](#): We employ a toolset --- dubbed Dr. Frankenstein --- to analyse the similarity of representations in deep neural networks. With this toolset we aim to match the activations on given layers of two trained neural networks by joining them with a stitching layer. We demonstrate that the inner representations emerging in deep convolutional neural networks with the same architecture but different initialisations can be matched with a surprisingly high degree of accuracy even with a single, affine stitching layer. We choose the stitching layer from several possible classes of linear transformations and investigate their performance and properties. The task of matching representations is closely related to notions of similarity. Using this toolset we also provide a novel viewpoint on the current line of research regarding similarity indices of neural network representations: the perspective of the performance on a task.

[DOCTOR: A Simple Method for Detecting Misclassification Errors](#)

- Federica Granese, Marco Romanelli, Daniele Gorla, Catuscia Palamidessi, Pablo Piantanida
- abstract@[open-review](#): Deep neural networks (DNNs) have shown to perform very well on large scale object recognition problems and lead to widespread use for real-world applications, including situations where DNN are implemented as “black boxes”. A promising approach to secure their use is to accept decisions that are likely to be correct while discarding the others. In this work, we propose DOCTOR, a simple method that aims to identify whether the prediction of a DNN classifier should (or should not) be trusted so that, consequently, it would be possible to accept it or to reject it. Two scenarios are investigated: Totally Black Box (TBB) where only the soft-predictions are available and Partially Black Box (PBB) where gradient-propagation to perform input pre-processing is allowed. Empirically, we show that DOCTOR outperforms all state-of-the-art methods on various well-known images and sentiment analysis datasets. In particular, we observe a reduction of up to 4% of the false rejection rate (FRR) in the PBB scenario. DOCTOR can be applied to any pre-trained model, it does not require prior information about the underlying dataset and is as simple as the simplest available methods in the literature.

[Contrastive Laplacian Eigenmaps](#)

- Hao Zhu, Ke Sun, Peter Koniusz
- abstract@[open-review](#): Graph contrastive learning attracts/disperses node representations for similar/dissimilar node pairs under some notion of similarity. It may be combined with a low-dimensional embedding of nodes to preserve intrinsic and structural properties of a graph. In this paper, we extend the celebrated Laplacian Eigenmaps with contrastive learning, and call them COntastive Laplacian EigenmapS (COLES). Starting from a GAN-inspired contrastive formulation, we show that the Jensen-Shannon divergence underlying many contrastive graph embedding models fails under disjoint positive and negative distributions, which may naturally emerge during sampling in the contrastive setting. In contrast, we demonstrate analytically that COLES essentially minimizes a surrogate of Wasserstein distance, which is known to cope well under disjoint distributions. Moreover, we show that the loss of COLES belongs to the family of so-called block-contrastive losses, previously shown to be superior compared to pair-wise losses typically used by contrastive methods. We show on popular benchmarks/backbones that COLES offers favourable accuracy/scalability compared to DeepWalk, GCN, Graph2Gauss, DGI and GRACE baselines.

[Machine learning structure preserving brackets for forecasting irreversible processes](#)

- Kookjin Lee, Nathaniel Trask, Panos Stinis
- abstract@[open-review](#): Forecasting of time-series data requires imposition of inductive biases to obtain predictive extrapolation, and recent works have imposed Hamiltonian/Lagrangian form to preserve structure for systems with \emph{reversible} dynamics. In this work we present a novel parameterization of dissipative brackets from metriplectic dynamical systems appropriate for learning \emph{irreversible} dynamics with unknown a priori model form. The process learns generalized Casimirs for energy and entropy guaranteed to be conserved and nondecreasing, respectively. Furthermore, for the case of added thermal noise, we guarantee exact preservation of a fluctuation-dissipation theorem, ensuring thermodynamic consistency. We provide benchmarks for dissipative systems demonstrating learned dynamics are more robust and generalize better than either "black-box" or penalty-based approaches.

[On the Variance of the Fisher Information for Deep Learning](#)

- Alexander Soen, Ke Sun
- abstract@[open-review](#): In the realm of deep learning, the Fisher information matrix (FIM) gives novel insights and useful tools to characterize the loss landscape, perform second-order optimization, and build geometric learning theories. The exact FIM is either unavailable in closed form or too expensive to compute. In practice, it is almost always estimated based on empirical samples. We investigate two such estimators based on two equivalent representations of the FIM --- both unbiased and consistent. Their estimation quality is naturally gauged by their variance given in closed form. We analyze how the parametric structure of a deep neural network can affect the variance. The meaning of this variance measure and its upper bounds are then discussed in the context of deep learning.

[A\\$^2\\$-Net: Learning Attribute-Aware Hash Codes for Large-Scale Fine-Grained Image Retrieval](#)

- Xiu-Shen Wei, Yang Shen, Xuhao Sun, Han-Jia Ye, Jian Yang
- abstract@[open-review](#): Our work focuses on tackling large-scale fine-grained image retrieval as ranking the images depicting the concept of interests (i.e., the same sub-category labels) highest based on the fine-grained details in the query. It is desirable to alleviate the challenges of both fine-grained nature of small inter-class variations with large intra-class variations and explosive growth of fine-grained data for such a practical task. In this paper, we propose an Attribute-Aware hashing Network (A\$^2\$-Net) for generating attribute-aware hash codes to not only make the retrieval process efficient, but also establish explicit correspondences between hash codes and visual attributes. Specifically, based on the captured visual representations by attention, we develop an encoder-decoder structure network of a reconstruction task to unsupervisedly distill high-level attribute-specific vectors from the appearance-specific visual representations without attribute annotations. A\$^2\$-Net is also equipped with a feature decorrelation constraint upon these attribute vectors to enhance their representation abilities. Finally, the required hash codes are generated by the attribute vectors driven by preserving original similarities. Qualitative experiments on five benchmark fine-grained datasets show our superiority over competing methods. More importantly, quantitative results demonstrate the obtained hash codes can strongly correspond to certain kinds of crucial properties of fine-grained objects.

[Shape Registration in the Time of Transformers](#)

- Giovanni Trappolini, Luca Cosmo, Luca Moschella, Riccardo Marin, Simone Melzi, Emanuele Rodolà
- abstract@[open-review](#): In this paper, we propose a transformer-based procedure for the efficient registration of non-rigid 3D point clouds. The proposed approach is data-driven and adopts for the first time the transformers architecture in the registration task. Our method is general and applies to different settings. Given a fixed template with some desired properties (e.g. skinning weights or other animation cues), we can register raw acquired data to it, thereby transferring all the template properties to the input geometry. Alternatively, given a pair of shapes, our method can register the first onto the second (or vice-versa), obtaining a high-quality dense correspondence between the two. In both contexts, the quality of our results enables us to target real applications such as texture transfer and shape interpolation. Furthermore, we also show that including an estimation of the underlying density of the surface eases the learning process. By exploiting the potential of this architecture, we can train our model requiring only a sparse set of ground truth correspondences (\$10\sim20\%\$ of the total points). The proposed model and the analysis that we perform pave the way for future exploration of

transformer-based architectures for registration and matching applications. Qualitative and quantitative evaluations demonstrate that our pipeline outperforms state-of-the-art methods for deformable and unordered 3D data registration on different datasets and scenarios.

[Brick-by-Brick: Combinatorial Construction with Deep Reinforcement Learning](#)

- Hyunsoo Chung, Jungtaek Kim, Boris Knyazev, Jinhwi Lee, Graham W. Taylor, Jaesik Park, Minsu Cho
- abstract@[open-review](#): Discovering a solution in a combinatorial space is prevalent in many real-world problems but it is also challenging due to diverse complex constraints and the vast number of possible combinations. To address such a problem, we introduce a novel formulation, combinatorial construction, which requires a building agent to assemble unit primitives (i.e., LEGO bricks) sequentially -- every connection between two bricks must follow a fixed rule, while no bricks mutually overlap. To construct a target object, we provide incomplete knowledge about the desired target (i.e., 2D images) instead of exact and explicit volumetric information to the agent. This problem requires a comprehensive understanding of partial information and long-term planning to append a brick sequentially, which leads us to employ reinforcement learning. The approach has to consider a variable-sized action space where a large number of invalid actions, which would cause overlap between bricks, exist. To resolve these issues, our model, dubbed Brick-by-Brick, adopts an action validity prediction network that efficiently filters invalid actions for an actor-critic network. We demonstrate that the proposed method successfully learns to construct an unseen object conditioned on a single image or multiple views of a target object.

[Dissecting the Diffusion Process in Linear Graph Convolutional Networks](#)

- Yifei Wang, Yisen Wang, Jiansheng Yang, Zhouchen Lin
- abstract@[open-review](#): Graph Convolutional Networks (GCNs) have attracted more and more attentions in recent years. A typical GCN layer consists of a linear feature propagation step and a nonlinear transformation step. Recent works show that a linear GCN can achieve comparable performance to the original non-linear GCN while being much more computationally efficient. In this paper, we dissect the feature propagation steps of linear GCNs from a perspective of continuous graph diffusion, and analyze why linear GCNs fail to benefit from more propagation steps. Following that, we propose Decoupled Graph Convolution (DGC) that decouples the terminal time and the feature propagation steps, making it more flexible and capable of exploiting a very large number of feature propagation steps. Experiments demonstrate that our proposed DGC improves linear GCNs by a large margin and makes them competitive with many modern variants of non-linear GCNs.

[Dynamic Grained Encoder for Vision Transformers](#)

- Lin Song, Songyang Zhang, Songtao Liu, Zeming Li, Xuming He, Hongbin Sun, Jian Sun, Nanning Zheng
- abstract@[open-review](#): Transformers, the de-facto standard for language modeling, have been recently applied for vision tasks. This paper introduces sparse queries for vision transformers to exploit the intrinsic spatial redundancy of natural images and save computational costs. Specifically, we propose a Dynamic Grained Encoder for vision transformers, which can adaptively assign a suitable number of queries to each spatial region. Thus it achieves a fine-grained representation in discriminative regions while keeping high efficiency. Besides, the dynamic grained encoder is compatible with most vision transformer frameworks. Without bells and whistles, our encoder allows the state-of-the-art vision transformers to reduce computational complexity by 40%-60% while maintaining comparable performance on image classification. Extensive experiments on object detection and segmentation further demonstrate the generalizability of our approach. Code is available at <https://github.com/StevenGrove/vtpack>.

[Understanding Negative Samples in Instance Discriminative Self-supervised Representation Learning](#)

- Kento Nozawa, Issei Sato
- abstract@[open-review](#): Instance discriminative self-supervised representation learning has been attracted attention thanks to its unsupervised nature and informative feature representation for downstream tasks. In practice, it commonly uses a larger number of negative samples than the number of supervised classes. However, there is an inconsistency in the existing analysis; theoretically, a large number of negative samples degrade classification performance on a downstream supervised task, while empirically, they improve the performance. We provide a novel framework to analyze this empirical result regarding negative samples using the coupon collector's problem. Our bound can implicitly incorporate the supervised loss of the downstream task in the self-supervised loss by increasing the number of negative samples. We confirm that our proposed analysis holds on real-world benchmark datasets.

[On UMAP's True Loss Function](#)

- Sebastian Damrich, Fred A. Hamprecht
- abstract@[open-review](#): UMAP has supplanted t-SNE as state-of-the-art for visualizing high-dimensional datasets in many disciplines, but the reason for its success is not well understood. In this work, we investigate UMAP's sampling based optimization scheme in detail. We derive UMAP's true loss function in closed form and find that it differs from the published one in a dataset size dependent way. As a consequence, we show that UMAP does not aim to reproduce its theoretically motivated high-dimensional UMAP similarities. Instead, it tries to reproduce similarities that only encode the k nearest neighbor graph, thereby challenging the previous understanding of UMAP's effectiveness. Alternatively, we consider the implicit balancing of attraction and repulsion due to the negative sampling to be key to UMAP's success. We corroborate our theoretical findings on toy and single cell RNA sequencing data.

[Fast Pure Exploration via Frank-Wolfe](#)

- Po-An Wang, Ruo-Chun Tzeng, Alexandre Proutiere
- abstract@[open-review](#): We study the problem of active pure exploration with fixed confidence in generic stochastic bandit environments. The goal of the learner is to answer a query about the environment with a given level of certainty while minimizing her sampling budget. For this problem, instance-specific lower bounds on the expected sample complexity reveal the optimal proportions of arm draws an Oracle algorithm would apply. These proportions solve an optimization problem whose tractability strongly depends on the structural properties of the environment, but may be instrumental in the design of efficient learning algorithms. We devise Frank-Wolfe-based Sampling (FWS), a simple algorithm whose sample complexity matches the lower bounds for a wide class of pure exploration problems. The algorithm is computationally efficient as, to learn and track the optimal proportion of arm draws, it relies on a single iteration of Frank-Wolfe algorithm applied to the lower-bound optimization problem. We apply FWS to various pure exploration tasks, including best arm identification in unstructured, thresholded, linear, and Lipschitz bandits. Despite its simplicity, FWS is competitive compared to state-of-art algorithms.

[iFlow: Numerically Invertible Flows for Efficient Lossless Compression via a Uniform Coder](#)

- Shifeng Zhang, Ning Kang, Tom Ryder, Zhenguo Li
- abstract@[open-review](#): It was estimated that the world produced 59 ZB (\$5.9 \times 10^{13}\$ GB) of data in 2020, resulting in the enormous costs of both data storage and transmission. Fortunately, recent advances in deep generative models have spearheaded a new class of so-called "neural compression" algorithms, which significantly outperform traditional codecs in terms of compression ratio. Unfortunately, the application of neural compression garners little commercial interest due to its limited bandwidth; therefore, developing highly efficient frameworks is of critical practical importance. In this paper, we discuss lossless compression using normalizing flows which have demonstrated a great capacity for achieving high compression ratios. As such, we introduce iFlow, a new method for achieving efficient lossless compression. We first propose Modular Scale Transform

(MST) and a novel family of numerically invertible flow transformations based on MST. Then we introduce the Uniform Base Conversion System (UBCS), a fast uniform-distribution codec incorporated into iFlow, enabling efficient compression. iFlow achieves state-of-the-art compression ratios and is \$5 \times\$ quicker than other high-performance schemes. Furthermore, the techniques presented in this paper can be used to accelerate coding time for a broad class of flow-based algorithms.

[History Aware Multimodal Transformer for Vision-and-Language Navigation](#)

- Shizhe Chen, Pierre-Louis Guhur, Cordelia Schmid, Ivan Laptev
- abstract@[open-review](#): Vision-and-language navigation (VLN) aims to build autonomous visual agents that follow instructions and navigate in real scenes. To remember previously visited locations and actions taken, most approaches to VLN implement memory using recurrent states. Instead, we introduce a History Aware Multimodal Transformer (HAMT) to incorporate a long-horizon history into multimodal decision making. HAMT efficiently encodes all the past panoramic observations via a hierarchical vision transformer (ViT), which first encodes individual images with ViT, then models spatial relation between images in a panoramic observation and finally takes into account temporal relation between panoramas in the history. It, then, jointly combines text, history and current observation to predict the next action. We first train HAMT end-to-end using several proxy tasks including single step action prediction and spatial relation prediction, and then use reinforcement learning to further improve the navigation policy. HAMT achieves new state of the art on a broad range of VLN tasks, including VLN with fine-grained instructions (R2R, RxR), high-level instructions (R2R-Last, REVERIE), dialogs (CVDN) as well as long-horizon VLN (R4R, R2R-Back). We demonstrate HAMT to be particularly effective for navigation tasks with longer trajectories.

[Meta Two-Sample Testing: Learning Kernels for Testing with Limited Data](#)

- Feng Liu, Wenkai Xu, Jie Lu, Danica J. Sutherland
- abstract@[open-review](#): Modern kernel-based two-sample tests have shown great success in distinguishing complex, high-dimensional distributions by learning appropriate kernels (or, as a special case, classifiers). Previous work, however, has assumed that many samples are observed from both of the distributions being distinguished. In realistic scenarios with very limited numbers of data samples, it can be challenging to identify a kernel powerful enough to distinguish complex distributions. We address this issue by introducing the problem of meta two-sample testing (M2ST), which aims to exploit (abundant) auxiliary data on related tasks to find an algorithm that can quickly identify a powerful test on new target tasks. We propose two specific algorithms for this task: a generic scheme which improves over baselines, and a more tailored approach which performs even better. We provide both theoretical justification and empirical evidence that our proposed meta-testing schemes outperform learning kernel-based tests directly from scarce observations, and identify when such schemes will be successful.

[Process for Adapting Language Models to Society \(PALMS\) with Values-Targeted Datasets](#)

- Irene Solaiman, Christy Dennison
- abstract@[open-review](#): Language models can generate harmful and biased outputs and exhibit undesirable behavior according to a given cultural context. We propose a Process for Adapting Language Models to Society (PALMS) with Values-Targeted Datasets, an iterative process to significantly change model behavior by crafting and fine-tuning on a dataset that reflects a predetermined set of target values. We evaluate our process using three metrics: quantitative metrics with human evaluations that score output adherence to a target value, toxicity scoring on outputs; and qualitative metrics analyzing the most common word associated with a given social category. Through each iteration, we add additional training dataset examples based on observed shortcomings from evaluations. PALMS performs significantly better on all metrics compared to baseline and control models for a broad range of GPT-3 language model sizes without compromising capability integrity. We find that the effectiveness of PALMS increases with model size. We show that significantly adjusting language model behavior is feasible with a small, hand-curated dataset.

[The Lazy Online Subgradient Algorithm is Universal on Strongly Convex Domains](#)

- Daron Anderson, Douglas Leith
- abstract@[open-review](#): We study Online Lazy Gradient Descent for optimisation on a strongly convex domain. The algorithm is known to achieve $\mathcal{O}(\sqrt{N})$ regret against adversarial opponents; here we show it is universal in the sense that it also achieves $\mathcal{O}(\log N)$ expected regret against i.i.d opponents. This improves upon the more complex meta-algorithm of Huang et al \cite{FTLBall} that only gets $\mathcal{O}(\sqrt{N \log N})$ and $\mathcal{O}(\log N)$ bounds. In addition we show that, unlike for the simplex, order bounds for pseudo-regret and expected regret are equivalent for strongly convex domains.

[Computer-Aided Design as Language](#)

- Yaroslav Ganin, Sergey Bartunov, Yujia Li, Ethan Keller, Stefano Saliceti
- abstract@[open-review](#): Computer-Aided Design (CAD) applications are used in manufacturing to model everything from coffee mugs to sports cars. These programs are complex and require years of training and experience to master. A component of all CAD models particularly difficult to make are the highly structured 2D sketches that lie at the heart of every 3D construction. In this work, we propose a machine learning model capable of automatically generating such sketches. Through this, we pave the way for developing intelligent tools that would help engineers create better designs with less effort. The core of our method is a combination of a general-purpose language modeling technique alongside an off-the-shelf data serialization protocol. Additionally, we explore several extensions allowing us to gain finer control over the generation process. We show that our approach has enough flexibility to accommodate the complexity of the domain and performs well for both unconditional synthesis and image-to-sketch translation.

[COHESIV: Contrastive Object and Hand Embedding Segmentation In Video](#)

- Dandan Shan, Richard Higgins, David Fouhey
- abstract@[open-review](#): In this paper we learn to segment hands and hand-held objects from motion. Our system takes a single RGB image and hand location as input to segment the hand and hand-held object. For learning, we generate responsibility maps that show how well a hand's motion explains other pixels' motion in video. We use these responsibility maps as pseudo-labels to train a weakly-supervised neural network using an attention-based similarity loss and contrastive loss. Our system outperforms alternate methods, achieving good performance on the 100DOH, EPIC-KITCHENS, and HO3D datasets.

[ByPE-VAE: Bayesian Pseudocoresets Exemplar VAE](#)

- Qingzhong Ai, LIRONG HE, SHIYU LIU, Zenglin Xu
- abstract@[open-review](#): Recent studies show that advanced priors play a major role in deep generative models. Exemplar VAE, as a variant of VAE with an exemplar-based prior, has achieved impressive results. However, due to the nature of model design, an exemplar-based model usually requires vast amounts of data to participate in training, which leads to huge computational complexity. To address this issue, we propose Bayesian Pseudocoresets Exemplar VAE (ByPE-VAE), a new variant of VAE with a prior based on Bayesian pseudocoreset. The proposed prior is conditioned on a small-scale pseudocoreset rather than the whole dataset for reducing the computational cost and avoiding overfitting. Simultaneously, we obtain the optimal pseudocoreset via a stochastic optimization algorithm during VAE training aiming to minimize the Kullback-Leibler divergence between the prior based on the pseudocoreset and that based on the whole dataset. Experimental results show that ByPE-VAE can achieve competitive improvements over the

state-of-the-art VAEs in the tasks of density estimation, representation learning, and generative data augmentation. Particularly, on a basic VAE architecture, ByPE-VAE is up to 3 times faster than Exemplar VAE while almost holding the performance. Code is available at \url{https://github.com/Aiqz/ByPE-VAE}.

[Recovery Analysis for Plug-and-Play Priors using the Restricted Eigenvalue Condition](#)

- Jiaming Liu, Salman Asif, Brendt Wohlberg, Ulugbek Kamilov
- abstract@[open-review](#): The plug-and-play priors (PnP) and regularization by denoising (RED) methods have become widely used for solving inverse problems by leveraging pre-trained deep denoisers as image priors. While the empirical imaging performance and the theoretical convergence properties of these algorithms have been widely investigated, their recovery properties have not previously been theoretically analyzed. We address this gap by showing how to establish theoretical recovery guarantees for PnP/RED by assuming that the solution of these methods lies near the fixed-points of a deep neural network. We also present numerical results comparing the recovery performance of PnP/RED in compressive sensing against that of recent compressive sensing algorithms based on generative models. Our numerical results suggest that PnP with a pre-trained artifact removal network provides significantly better results compared to the existing state-of-the-art methods.

[Group Equivariant Subsampling](#)

- Jin Xu, Hyunjik Kim, Thomas Rainforth, Yee Teh
- abstract@[open-review](#): Subsampling is used in convolutional neural networks (CNNs) in the form of pooling or strided convolutions, to reduce the spatial dimensions of feature maps and to allow the receptive fields to grow exponentially with depth. However, it is known that such subsampling operations are not translation equivariant, unlike convolutions that are translation equivariant. Here, we first introduce translation equivariant subsampling/upsampling layers that can be used to construct exact translation equivariant CNNs. We then generalise these layers beyond translations to general groups, thus proposing group equivariant subsampling/upsampling. We use these layers to construct group equivariant autoencoders (GAEs) that allow us to learn low-dimensional equivariant representations. We empirically verify on images that the representations are indeed equivariant to input translations and rotations, and thus generalise well to unseen positions and orientations. We further use GAEs in models that learn object-centric representations on multi-object datasets, and show improved data efficiency and decomposition compared to non-equivariant baselines.

[Data Sharing and Compression for Cooperative Networked Control](#)

- Jiangan Cheng, Marco Pavone, Sachin Katti, Sandeep Chinchali, Ao Tang
- abstract@[open-review](#): Sharing forecasts of network timeseries data, such as cellular or electricity load patterns, can improve independent control applications ranging from traffic scheduling to power generation. Typically, forecasts are designed without knowledge of a downstream controller's task objective, and thus simply optimize for mean prediction error. However, such task-agnostic representations are often too large to stream over a communication network and do not emphasize salient temporal features for cooperative control. This paper presents a solution to learn succinct, highly-compressed forecasts that are co-designed with a modular controller's task objective. Our simulations with real cellular, Internet-of-Things (IoT), and electricity load data show we can improve a model predictive controller's performance by at least 25% while transmitting 80% less data than the competing method. Further, we present theoretical compression results for a networked variant of the classical linear quadratic regulator (LQR) control problem.

[Hyperbolic Procrustes Analysis Using Riemannian Geometry](#)

- Ya-Wei Eileen Lin, Yuval Kluger, Ronen Talmon
- abstract@[open-review](#): Label-free alignment between datasets collected at different times, locations, or by different instruments is a fundamental scientific task. Hyperbolic spaces have recently provided a fruitful foundation for the development of informative representations of hierarchical data. Here, we take a purely geometric approach for label-free alignment of hierarchical datasets and introduce hyperbolic Procrustes analysis (HPA). HPA consists of new implementations of the three prototypical Procrustes analysis components: translation, scaling, and rotation, based on the Riemannian geometry of the Lorentz model of hyperbolic space. We analyze the proposed components, highlighting their useful properties for alignment. The efficacy of HPA, its theoretical properties, stability and computational efficiency are demonstrated in simulations. In addition, we showcase its performance on three batch correction tasks involving gene expression and mass cytometry data. Specifically, we demonstrate high-quality unsupervised batch effect removal from data acquired at different sites and with different technologies that outperforms recent methods for label-free alignment in hyperbolic spaces.

[No Fear of Heterogeneity: Classifier Calibration for Federated Learning with Non-IID Data](#)

- Mi Luo, Fei Chen, Dapeng Hu, Yifan Zhang, Jian Liang, Jiashi Feng
- abstract@[open-review](#): A central challenge in training classification models in the real-world federated system is learning with non-IID data. To cope with this, most of the existing works involve enforcing regularization in local optimization or improving the model aggregation scheme at the server. Other works also share public datasets or synthesized samples to supplement the training of under-represented classes or introduce a certain level of personalization. Though effective, they lack a deep understanding of how the data heterogeneity affects each layer of a deep classification model. In this paper, we bridge this gap by performing an experimental analysis of the representations learned by different layers. Our observations are surprising: (1) there exists a greater bias in the classifier than other layers, and (2) the classification performance can be significantly improved by post-calibrating the classifier after federated training. Motivated by the above findings, we propose a novel and simple algorithm called Classifier Calibration with Virtual Representations (CCVR), which adjusts the classifier using virtual representations sampled from an approximated gaussian mixture model. Experimental results demonstrate that CCVR achieves state-of-the-art performance on popular federated learning benchmarks including CIFAR-10, CIFAR-100, and CINIC-10. We hope that our simple yet effective method can shed some light on the future research of federated learning with non-IID data.

[Preconditioned Gradient Descent for Over-Parameterized Nonconvex Matrix Factorization](#)

- Jialun Zhang, Salar Fattah, Richard Y Zhang
- abstract@[open-review](#): In practical instances of nonconvex matrix factorization, the rank of the true solution $\$r^{\star}\$$ is often unknown, so the rank $\$r\$$ of the model can be over-specified as $\$r > r^{\star}\$$. This over-parameterized regime of matrix factorization significantly slows down the convergence of local search algorithms, from a linear rate with $\$r = r^{\star}\$$ to a sublinear rate when $\$r > r^{\star}\$$. We propose an inexpensive preconditioner for the matrix sensing variant of nonconvex matrix factorization that restores the convergence rate of gradient descent back to linear, even in the over-parameterized case, while also making it agnostic to possible ill-conditioning in the ground truth. Classical gradient descent in a neighborhood of the solution slows down due to the need for the model matrix factor to become singular. Our key result is that this singularity can be corrected by $\$ell_2\$$ regularization with a specific range of values for the damping parameter. In fact, a good damping parameter can be inexpensively estimated from the current iterate. The resulting algorithm, which we call preconditioned gradient descent or PrecGD, is stable under noise, and converges linearly to an information theoretically optimal error bound. Our numerical experiments find that PrecGD works equally well in restoring the linear convergence of other variants of nonconvex matrix factorization in the over-parameterized regime.

[Improving Contrastive Learning on Imbalanced Data via Open-World Sampling](#)

- Ziyu Jiang, Tianlong Chen, Ting Chen, Zhangyang Wang
- abstract@[open-review](#): Contrastive learning approaches have achieved great success in learning visual representations with few labels of the target classes. That implies a tantalizing possibility of scaling them up beyond a curated “seed” benchmark, to incorporating more unlabeled images from the internet-scale external sources to enhance its performance. However, in practice, larger amount of unlabeled data will require more computing resources due to the bigger model size and longer training needed. Moreover, open-world unlabeled data usually follows an implicit long-tail class or attribute distribution, many of which also do not belong to the target classes. Blindly leveraging all unlabeled data hence can lead to the data imbalance as well as distraction issues. This motivates us to seek a principled approach to strategically select unlabeled data from an external source, in order to learn generalizable, balanced and diverse representations for relevant classes. In this work, we present an open-world unlabeled data sampling framework called Model-Aware K-center (MAK), which follows three simple principles: (1) tailness, which encourages sampling of examples from tail classes, by sorting the empirical contrastive loss expectation (ECLE) of samples over random data augmentations; (2) proximity, which rejects the out-of-distribution outliers that may distract training; and (3) diversity, which ensures diversity in the set of sampled examples. Empirically, using ImageNet-100-LT (without labels) as the seed dataset and two “noisy” external data sources, we demonstrate that MAK can consistently improve both the overall representation quality and the class balancedness of the learned features, as evaluated via linear classifier evaluation on full-shot and few-shot settings. The code is available at: <https://github.com/VITA-Group/MAK>.

[Searching for Efficient Transformers for Language Modeling](#)

- David So, Wojciech Mańke, Hanxiao Liu, Zihang Dai, Noam Shazeer, Quoc V Le
- abstract@[open-review](#): Large Transformer models have been central to recent advances in natural language processing. The training and inference costs of these models, however, have grown rapidly and become prohibitively expensive. Here we aim to reduce the costs of Transformers by searching for a more efficient variant. Compared to previous approaches, our search is performed at a lower level, over the primitives that define a Transformer TensorFlow program. We identify an architecture, named Primer, that has a smaller training cost than the original Transformer and other variants for auto-regressive language modeling. Primer’s improvements can be mostly attributed to two simple modifications: squaring ReLU activations and adding a depthwise convolution layer after each Q, K, and V projection in self-attention. Experiments show Primer’s gains over Transformer increase as compute scale grows and follow a power law with respect to quality at optimal model sizes. We also verify empirically that Primer can be dropped into different codebases to significantly speed up training without additional tuning. For example, at a 500M parameter size, Primer improves the original T5 architecture on C4 auto-regressive language modeling, reducing the training cost by 4X. Furthermore, the reduced training cost means Primer needs much less compute to reach a target one-shot performance. For instance, in a 1.9B parameter configuration similar to GPT-3 XL, Primer uses 1/3 of the training compute to achieve the same one-shot performance as Transformer. We open source our models and several comparisons in T5 to help with reproducibility.

[Scaling Ensemble Distribution Distillation to Many Classes with Proxy Targets](#)

- Max Ryabinin, Andrey Malinin, Mark Gales
- abstract@[open-review](#): Ensembles of machine learning models yield improved system performance as well as robust and interpretable uncertainty estimates; however, their inference costs can be prohibitively high. Ensemble Distribution Distillation (EnD\$^2\$) is an approach that allows a single model to efficiently capture both the predictive performance and uncertainty estimates of an ensemble. For classification, this is achieved by training a Dirichlet distribution over the ensemble members’ output distributions via the maximum likelihood criterion. Although theoretically principled, this work shows that the criterion exhibits poor convergence when applied to large-scale tasks where the number of classes is very high. Specifically, we show that for the Dirichlet log-likelihood criterion classes with low probability induce larger gradients than high-probability classes. Hence during training the model focuses on the distribution of the ensemble tail-class probabilities rather than the probability of the correct and closely related classes. We propose a new training objective which minimizes the reverse KL-divergence to a \emph{Proxy-Dirichlet} target derived from the ensemble. This loss resolves the gradient issues of EnD\$^2\$, as we demonstrate both theoretically and empirically on the ImageNet, LibriSpeech, and WMT17 En-De datasets containing 1000, 5000, and 40,000 classes, respectively.

[Multi-Person 3D Motion Prediction with Multi-Range Transformers](#)

- Jiashun Wang, Huazhe Xu, Medhini Narasimhan, Xiaolong Wang
- abstract@[open-review](#): We propose a novel framework for multi-person 3D motion trajectory prediction. Our key observation is that a human’s action and behaviors may highly depend on the other persons around. Thus, instead of predicting each human pose trajectory in isolation, we introduce a Multi-Range Transformers model which contains of a local-range encoder for individual motion and a global-range encoder for social interactions. The Transformer decoder then performs prediction for each person by taking a corresponding pose as a query which attends to both local and global-range encoder features. Our model not only outperforms state-of-the-art methods on long-term 3D motion prediction, but also generates diverse social interactions. More interestingly, our model can even predict 15-person motion simultaneously by automatically dividing the persons into different interaction groups. Project page with code is available at <https://jiashunwang.github.io/MRT/>.

[STEM: A Stochastic Two-Sided Momentum Algorithm Achieving Near-Optimal Sample and Communication Complexities for Federated Learning](#)

- Prashant Khanduri, PRANAY SHARMA, Haibo Yang, Mingyi Hong, Jia Liu, Ketan Rajawat, Pramod Varshney
- abstract@[open-review](#): Federated Learning (FL) refers to the paradigm where multiple worker nodes (WNs) build a joint model by using local data. Despite extensive research, for a generic non-convex FL problem, it is not clear, how to choose the WNs’ and the server’s update directions, the minibatch sizes, and the local update frequency, so that the WNs use the minimum number of samples and communication rounds to achieve the desired solution. This work addresses the above question and considers a class of stochastic algorithms where the WNs perform a few local updates before communication. We show that when both the WN’s and the server’s directions are chosen based on certain stochastic momentum estimator, the algorithm requires $\tilde{\mathcal{O}}(\epsilon^{-3/2})$ samples and $\tilde{\mathcal{O}}(\epsilon^{-1})$ communication rounds to compute an ϵ -stationary solution. To the best of our knowledge, this is the first FL algorithm that achieves such \{it near-optimal\} sample and communication complexities simultaneously. Further, we show that there is a trade-off curve between local update frequencies and local minibatch sizes, on which the above sample and communication complexities can be maintained. Finally, we show that for the classical FedAvg (a.k.a. Local SGD, which is a momentum-less special case of the STEM), a similar trade-off curve exists, albeit with worse sample and communication complexities. Our insights on this trade-off provides guidelines for choosing the four important design elements for FL algorithms, the update frequency, directions, and minibatch sizes to achieve the best performance.}

[Bubblewrap: Online tiling and real-time flow prediction on neural manifolds](#)

- Anne Draelos, Pranjali Gupta, Na Young Jun, Chaichontat Sriworarat, John Pearson
- abstract@[open-review](#): While most classic studies of function in experimental neuroscience have focused on the coding properties of individual neurons, recent developments in recording technologies have resulted in an increasing emphasis on the dynamics of neural populations. This has given rise to a wide variety of models for analyzing population activity in relation to experimental variables, but direct testing of many neural population hypotheses requires intervening in the system based on current neural state, necessitating models capable of inferring neural state online. Existing approaches, primarily based on dynamical systems, require strong parametric assumptions that are easily violated in the noise-dominated regime and do not scale well to the thousands of data channels in modern experiments. To address this problem, we propose a method that combines fast, stable dimensionality reduction with a soft tiling of the resulting neural manifold, allowing dynamics to be approximated as a probability flow between tiles. This method can be

fit efficiently using online expectation maximization, scales to tens of thousands of tiles, and outperforms existing methods when dynamics are noise-dominated or feature multi-modal transition probabilities. The resulting model can be trained at kiloHertz data rates, produces accurate approximations of neural dynamics within minutes, and generates predictions on submillisecond time scales. It retains predictive performance throughout many time steps into the future and is fast enough to serve as a component of closed-loop causal experiments.

[The Semi-Random Satisfaction of Voting Axioms](#)

- Lirong Xia
- abstract@[open-review](#): We initiate the work towards a comprehensive picture of the worst average-case satisfaction of voting axioms in semi-random models, to provide a finer and more realistic foundation for comparing voting rules. We adopt the semi-random model and formulation in [Xia 2020], where an adversary chooses arbitrarily correlated ``ground truth'' preferences for the agents, on top of which random noises are added. We focus on characterizing the semi-random satisfaction of two well-studied voting axioms: Condorcet criterion and participation. We prove that for any fixed number of alternatives, when the number of voters n is sufficiently large, the semi-random satisfaction of the Condorcet criterion under a wide range of voting rules is $\$1, \$1 \cdot \exp(-\sqrt{\Theta(n)})$, $\sqrt{\Theta(n^{0.5})}$, $\sqrt{\exp(-\sqrt{\Theta(n)})}$, or being $\sqrt{\Theta(1)}$ and $\sqrt{1 - \Theta(1)}$ at the same time; and the semi-random satisfaction of participation is $\sqrt{1 - \Theta(n^{0.5})}$. Our results address open questions by Berg and Lepelley in 1994, and also confirm the following high-level message: the Condorcet criterion is a bigger concern than participation under realistic models.

[Deep Marching Tetrahedra: a Hybrid Representation for High-Resolution 3D Shape Synthesis](#)

- Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, Sanja Fidler
- abstract@[open-review](#): We introduce DMTet, a deep 3D conditional generative model that can synthesize high-resolution 3D shapes using simple user guides such as coarse voxels. It marries the merits of implicit and explicit 3D representations by leveraging a novel hybrid 3D representation. Compared to the current implicit approaches, which are trained to regress the signed distance values, DMTet directly optimizes for the reconstructed surface, which enables us to synthesize finer geometric details with fewer artifacts. Unlike deep 3D generative models that directly generate explicit representations such as meshes, our model can synthesize shapes with arbitrary topology. The core of DMTet includes a deformable tetrahedral grid that encodes a discretized signed distance function and a differentiable marching tetrahedra layer that converts the implicit signed distance representation to the explicit surface mesh representation. This combination allows joint optimization of the surface geometry and topology as well as generation of the hierarchy of subdivisions using reconstruction and adversarial losses defined explicitly on the surface mesh. Our approach significantly outperforms existing work on conditional shape synthesis from coarse voxel inputs, trained on a dataset of complex 3D animal shapes. Project page: <https://nv-tlabs.github.io/DMTet/>.

[Learning to Combine Per-Example Solutions for Neural Program Synthesis](#)

- Disha Srivastava, Hugo Larochelle, Daniel Tarlow
- abstract@[open-review](#): The goal of program synthesis from examples is to find a computer program that is consistent with a given set of input-output examples. Most learning-based approaches try to find a program that satisfies all examples at once. Our work, by contrast, considers an approach that breaks the problem into two stages: (a) find programs that satisfy only one example, and (b) leverage these per-example solutions to yield a program that satisfies all examples. We introduce the Cross Aggregator neural network module based on a multi-head attention mechanism that learns to combine the cues present in these per-example solutions to synthesize a global solution. Evaluation across programs of different lengths and under two different experimental settings reveal that when given the same time budget, our technique significantly improves the success rate over PCCoder [Zohar et. al 2018] and other ablation baselines.

[On Success and Simplicity: A Second Look at Transferable Targeted Attacks](#)

- Zhengyu Zhao, Zhuoran Liu, Martha Larson
- abstract@[open-review](#): Achieving transferability of targeted attacks is reputed to be remarkably difficult. The current state of the art has resorted to resource-intensive solutions that necessitate training model(s) for each target class with additional data. In our investigation, we find, however, that simple transferable attacks which require neither model training nor additional data can achieve surprisingly strong targeted transferability. This insight has been overlooked until now, mainly because the widespread practice of attacking with only few iterations has largely limited the attack convergence to optimal targeted transferability. In particular, we, for the first time, identify that a very simple logit loss can largely surpass the commonly adopted cross-entropy loss, and yield even better results than the resource-intensive state of the art. Our analysis spans a variety of transfer scenarios, especially including three new, realistic scenarios: an ensemble transfer scenario with little model similarity, a worse-case scenario with low-ranked target classes, and also a real-world attack on the Google Cloud Vision API. Results in these new transfer scenarios demonstrate that the commonly adopted, easy scenarios cannot fully reveal the actual strength of different attacks and may cause misleading comparative results. We also show the usefulness of the simple logit loss for generating targeted universal adversarial perturbations in a data-free manner. Overall, the aim of our analysis is to inspire a more meaningful evaluation on targeted transferability. Code is available at <https://github.com/ZhengyuZhao/Targeted-Transfer>.

[Provably efficient, succinct, and precise explanations](#)

- Guy Blanc, Jane Lange, Li-Yang Tan
- abstract@[open-review](#): We consider the problem of explaining the predictions of an arbitrary blackbox model f : given query access to f and an instance x , output a small set of x 's features that in conjunction essentially determines $f(x)$. We design an efficient algorithm with provable guarantees on the succinctness and precision of the explanations that it returns. Prior algorithms were either efficient but lacked such guarantees, or achieved such guarantees but were inefficient. We obtain our algorithm via a connection to the problem of {\it implicitly} learning decision trees. The implicit nature of this learning task allows for efficient algorithms even when the complexity of f necessitates an intractably large surrogate decision tree. We solve the implicit learning problem by bringing together techniques from learning theory, local computation algorithms, and complexity theory. Our approach of “explaining by implicit learning” shares elements of two previously disparate methods for post-hoc explanations, global and local explanations, and we make the case that it enjoys advantages of both.

[Refined Learning Bounds for Kernel and Approximate \$k\$ -Means](#)

- Yong Liu
- abstract@[open-review](#): Kernel k -means is one of the most popular approaches to clustering and its theoretical properties have been investigated for decades. However, the existing state-of-the-art risk bounds are of order $\mathcal{O}(k\sqrt{n})$, which do not match with the stated lower bound $\Omega(\sqrt{k/n})$ in terms of k , where k is the number of clusters and n is the size of the training set. In this paper, we study the statistical properties of kernel k -means and Nyström-based kernel k -means, and obtain optimal clustering risk bounds, which improve the existing risk bounds. Particularly, based on a refined upper bound of Rademacher complexity [21], we first derive an optimal risk bound of rate $\mathcal{O}(\sqrt{k/n})$ for empirical risk minimizer (ERM), and further extend it to general cases beyond ERM. Then, we analyze the statistical effect of computational approximations of Nyström-based kernel k -means, and prove that it achieves the same statistical accuracy as the original kernel k -means considering only $\Omega(\sqrt{nk})$ Nyström landmark points. We further relax the restriction of landmark points from $\Omega(\sqrt{nk})$ to $\Omega(\sqrt{n})$ under a mild condition. Finally, we validate the theoretical findings via numerical experiments.

[Learning Causal Semantic Representation for Out-of-Distribution Prediction](#)

- Chang Liu, Xinwei Sun, Jindong Wang, Haoyue Tang, Tao Li, Tao Qin, Wei Chen, Tie-Yan Liu
- abstract@[open-review](#): Conventional supervised learning methods, especially deep ones, are found to be sensitive to out-of-distribution (OOD) examples, largely because the learned representation mixes the semantic factor with the variation factor due to their domain-specific correlation, while only the semantic factor causes the output. To address the problem, we propose a Causal Semantic Generative model (CSG) based on a causal reasoning so that the two factors are modeled separately, and develop methods for OOD prediction from a single training domain, which is common and challenging. The methods are based on the causal invariance principle, with a novel design in variational Bayes for both efficient learning and easy prediction. Theoretically, we prove that under certain conditions, CSG can identify the semantic factor by fitting training data, and this semantic-identification guarantees the boundedness of OOD generalization error and the success of adaptation. Empirical study shows improved OOD performance over prevailing baselines.

[A first-order primal-dual method with adaptivity to local smoothness](#)

- Maria-Luiza Vladarean, Yura Malitsky, Volkan Cevher
- abstract@[open-review](#): We consider the problem of finding a saddle point for the convex-concave objective $\min_x \max_y f(x) + \langle Ax, y \rangle - g^*(y)$, where f is a convex function with locally Lipschitz gradient and g is convex and possibly non-smooth. We propose an adaptive version of the Condat-Vũ algorithm, which alternates between primal gradient steps and dual proximal steps. The method achieves stepsize adaptivity through a simple rule involving $\|A\|$ and the norm of recently computed gradients of f . Under standard assumptions, we prove an $\mathcal{O}(k^{-1})$ ergodic convergence rate. Furthermore, when f is also locally strongly convex and A has full row rank we show that our method converges with a linear rate. Numerical experiments are provided for illustrating the practical performance of the algorithm.

[A Theory-Driven Self-Labeling Refinement Method for Contrastive Representation Learning](#)

- Pan Zhou, Caiming Xiong, Xiaotong Yuan, Steven Chu Hong Hoi
- abstract@[open-review](#): For an image query, unsupervised contrastive learning labels crops of the same image as positives, and other image crops as negatives. Although intuitive, such a native label assignment strategy cannot reveal the underlying semantic similarity between a query and its positives and negatives, and impairs performance, since some negatives are semantically similar to the query or even share the same semantic class as the query. In this work, we first prove that for contrastive learning, inaccurate label assignment heavily impairs its generalization for semantic instance discrimination, while accurate labels benefit its generalization. Inspired by this theory, we propose a novel self-labeling refinement approach for contrastive learning. It improves the label quality via two complementary modules: (i) self-labeling refinery (SLR) to generate accurate labels and (ii) momentum mixup (MM) to enhance similarity between query and its positive. SLR uses a positive of a query to estimate semantic similarity between a query and its positive and negatives, and combines estimated similarity with vanilla label assignment in contrastive learning to iteratively generate more accurate and informative soft labels. We theoretically show that our SLR can exactly recover the true semantic labels of label-corrupted data, and supervises networks to achieve zero prediction error on classification tasks. MM randomly combines queries and positives to increase semantic similarity between the generated virtual queries and their positives so as to improve label accuracy. Experimental results on CIFAR10, ImageNet, VOC and COCO show the effectiveness of our method.

[Adversarial Robustness with Semi-Infinite Constrained Learning](#)

- Alexander Robey, Luiz Chamon, George J. Pappas, Hamed Hassani, Alejandro Ribeiro
- abstract@[open-review](#): Despite strong performance in numerous applications, the fragility of deep learning to input perturbations has raised serious questions about its use in safety-critical domains. While adversarial training can mitigate this issue in practice, state-of-the-art methods are increasingly application-dependent, heuristic in nature, and suffer from fundamental trade-offs between nominal performance and robustness. Moreover, the problem of finding worst-case perturbations is non-convex and underparameterized, both of which engender a non-favorable optimization landscape. Thus, there is a gap between the theory and practice of robust learning, particularly with respect to when and why adversarial training works. In this paper, we take a constrained learning approach to address these questions and to provide a theoretical foundation for robust learning. In particular, we leverage semi-infinite optimization and non-convex duality theory to show that adversarial training is equivalent to a statistical problem over perturbation distributions. Notably, we show that a myriad of previous robust training techniques can be recovered for particular, sub-optimal choices of these distributions. Using these insights, we then propose a hybrid Langevin Markov Chain Monte Carlo approach for which several common algorithms (e.g., PGD) are special cases. Finally, we show that our approach can mitigate the trade-off between nominal and robust performance, yielding state-of-the-art results on MNIST and CIFAR-10. Our code is available at: <https://github.com/arobey1/advbench>.

[Conformal Time-series Forecasting](#)

- Kamile Stankeviciute, Ahmed M. Alaa, Mihaela van der Schaar
- abstract@[open-review](#): Current approaches for multi-horizon time series forecasting using recurrent neural networks (RNNs) focus on issuing point estimates, which is insufficient for decision-making in critical application domains where an uncertainty estimate is also required. Existing approaches for uncertainty quantification in RNN-based time-series forecasts are limited as they may require significant alterations to the underlying model architecture, may be computationally complex, may be difficult to calibrate, may incur high sample complexity, and may not provide theoretical guarantees on frequentist coverage. In this paper, we extend the inductive conformal prediction framework to the time-series forecasting setup, and propose a lightweight algorithm to address all of the above limitations, providing uncertainty estimates with theoretical guarantees for any multi-horizon forecast predictor and any dataset with minimal exchangeability assumptions. We demonstrate the effectiveness of our approach by comparing it with existing benchmarks on a variety of synthetic and real-world datasets.

[A 3D Generative Model for Structure-Based Drug Design](#)

- Shitong Luo, Jiaqi Guan, Jianzhu Ma, Jian Peng
- abstract@[open-review](#): We study a fundamental problem in structure-based drug design --- generating molecules that bind to specific protein binding sites. While we have witnessed the great success of deep generative models in drug design, the existing methods are mostly string-based or graph-based. They are limited by the lack of spatial information and thus unable to be applied to structure-based design tasks. Particularly, such models have no or little knowledge of how molecules interact with their target proteins exactly in 3D space. In this paper, we propose a 3D generative model that generates molecules given a designated 3D protein binding site. Specifically, given a binding site as the 3D context, our model estimates the probability density of atom's occurrences in 3D space --- positions that are more likely to have atoms will be assigned higher probability. To generate 3D molecules, we propose an auto-regressive sampling scheme --- atoms are sampled sequentially from the learned distribution until there is no room for new atoms. Combined with this sampling scheme, our model can generate valid and diverse molecules, which could be applicable to various structure-based molecular design tasks such as molecule sampling and linker design. Experimental results demonstrate that molecules sampled from our model exhibit high binding affinity to specific targets and good drug properties such as drug-likeness even if the model is not explicitly optimized for them.

[Bootstrapping the Error of Oja's Algorithm](#)

- Robert Lunde, Purnamrita Sarkar, Rachel Ward
- abstract@[open-review](#): We consider the problem of quantifying uncertainty for the estimation error of the leading eigenvector from Oja's algorithm for streaming principal component analysis, where the data are generated IID from some unknown distribution. By combining classical tools from the U-statistics literature with recent results on high-dimensional central limit theorems for quadratic forms of random vectors and concentration of matrix products, we establish a weighted χ^2 approximation result for the \sin^2 error between the population eigenvector and the output of Oja's algorithm. Since estimating the covariance matrix associated with the approximating distribution requires knowledge of unknown model parameters, we propose a multiplier bootstrap algorithm that may be updated in an online manner. We establish conditions under which the bootstrap distribution is close to the corresponding sampling distribution with high probability, thereby establishing the bootstrap as a consistent inferential method in an appropriate asymptotic regime.

[Landscape analysis of an improved power method for tensor decomposition](#)

- Joe Kileel, Timo Klock, João M Pereira
- abstract@[open-review](#): In this work, we consider the optimization formulation for symmetric tensor decomposition recently introduced in the Subspace Power Method (SPM) of Kileel and Pereira. Unlike popular alternative functionals for tensor decomposition, the SPM objective function has the desirable properties that its maximal value is known in advance, and its global optima are exactly the rank-1 components of the tensor when the input is sufficiently low-rank. We analyze the non-convex optimization landscape associated with the SPM objective. Our analysis accounts for working with noisy tensors. We derive quantitative bounds such that any second-order critical point with SPM objective value exceeding the bound must equal a tensor component in the noiseless case, and must approximate a tensor component in the noisy case. For decomposing tensors of size $D^{\times m}$, we obtain a near-global guarantee up to rank $\widetilde{o}(D^{\lfloor m/2 \rfloor})$ under a random tensor model, and a global guarantee up to rank $\mathcal{O}(D)$ assuming deterministic frame conditions. This implies that SPM with suitable initialization is a provable, efficient, robust algorithm for low-rank symmetric tensor decomposition. We conclude with numerics that show a practical preferability for using the SPM functional over a more established counterpart.

[Curriculum Offline Imitating Learning](#)

- Minghuan Liu, Hanye Zhao, Zhengyu Yang, Jian Shen, Weinan Zhang, Li Zhao, Tie-Yan Liu
- abstract@[open-review](#): Offline reinforcement learning (RL) tasks require the agent to learn from a pre-collected dataset with no further interactions with the environment. Despite the potential to surpass the behavioral policies, RL-based methods are generally impractical due to the training instability and bootstrapping the extrapolation errors, which always require careful hyperparameter tuning via online evaluation. In contrast, offline imitation learning (IL) has no such issues since it learns the policy directly without estimating the value function by bootstrapping. However, IL is usually limited in the capability of the behavioral policy and tends to learn a mediocre behavior from the dataset collected by the mixture of policies. In this paper, we aim to take advantage of IL but mitigate such a drawback. Observing that behavior cloning is able to imitate neighboring policies with less data, we propose Curriculum Offline Imitation Learning (COIL), which utilizes an experience picking strategy to make the agent imitate from adaptive neighboring policies with a higher return, and improves the current policy along curriculum stages. On continuous control benchmarks, we compare COIL against both imitation-based methods and RL-based methods, showing that COIL not only avoids just learning a mediocre behavior on mixed datasets but is also even competitive with state-of-the-art offline RL methods.

[Robust Pose Estimation in Crowded Scenes with Direct Pose-Level Inference](#)

- Dongkai Wang, Shiliang Zhang, Gang Hua
- abstract@[open-review](#): Multi-person pose estimation in crowded scenes is challenging because overlapping and occlusions make it difficult to detect person bounding boxes and infer pose cues from individual keypoints. To address those issues, this paper proposes a direct pose-level inference strategy that is free of bounding box detection and keypoint grouping. Instead of inferring individual keypoints, the Pose-level Inference Network (PINet) directly infers the complete pose cues for a person from his/her visible body parts. PINet first applies the Part-based Pose Generation (PPG) to infer multiple coarse poses for each person from his/her body parts. Those coarse poses are refined by the Pose Refinement module through incorporating pose priors, and finally are fused in the Pose Fusion module. PINet relies on discriminative body parts to differentiate overlapped persons, and applies visual body cues to infer the global pose cues. Experiments on several crowded scenes pose estimation benchmarks demonstrate the superiority of PINet. For instance, it achieves 59.8% AP on the OCHuman dataset, outperforming the recent works by a large margin.

[Ising Model Selection Using \$\ell_1\$ -Regularized Linear Regression: A Statistical Mechanics Analysis](#)

- Xiangming Meng, Tomoyuki Obuchi, Yoshiyuki Kabashima
- abstract@[open-review](#): We theoretically analyze the typical learning performance of ℓ_1 -regularized linear regression (ℓ_1 -LinR) for Ising model selection using the replica method from statistical mechanics. For typical random regular graphs in the paramagnetic phase, an accurate estimate of the typical sample complexity of ℓ_1 -LinR is obtained. Remarkably, despite the model misspecification, ℓ_1 -LinR is model selection consistent with the same order of sample complexity as ℓ_1 -regularized logistic regression (ℓ_1 -LogR), i.e., $M = \mathcal{O}(\log N)$, where N is the number of variables of the Ising model. Moreover, we provide an efficient method to accurately predict the non-asymptotic behavior of ℓ_1 -LinR for moderate M, N , such as precision and recall. Simulations show a fairly good agreement between theoretical predictions and experimental results, even for graphs with many loops, which supports our findings. Although this paper mainly focuses on ℓ_1 -LinR, our method is readily applicable for precisely characterizing the typical learning performances of a wide class of ℓ_1 -regularized M -estimators including ℓ_1 -LogR and interaction screening.

[Conformal Prediction using Conditional Histograms](#)

- Matteo Sesia, Yaniv Romano
- abstract@[open-review](#): This paper develops a conformal method to compute prediction intervals for non-parametric regression that can automatically adapt to skewed data. Leveraging black-box machine learning algorithms to estimate the conditional distribution of the outcome using histograms, it translates their output into the shortest prediction intervals with approximate conditional coverage. The resulting prediction intervals provably have marginal coverage in finite samples, while asymptotically achieving conditional coverage and optimal length if the black-box model is consistent. Numerical experiments with simulated and real data demonstrate improved performance compared to state-of-the-art alternatives, including conformalized quantile regression and other distributional conformal prediction approaches.

[Contrastive Graph Poisson Networks: Semi-Supervised Learning with Extremely Limited Labels](#)

- Sheng Wan, Yibing Zhan, Liu Liu, Baosheng Yu, Shirui Pan, Chen Gong
- abstract@[open-review](#): Graph Neural Networks (GNNs) have achieved remarkable performance in the task of semi-supervised node classification. However, most existing GNN models require sufficient labeled data for effective network training. Their performance can be seriously degraded when labels are extremely limited. To address this issue, we propose a new framework termed Contrastive Graph Poisson Networks (CGPN) for node classification under extremely limited labeled data. Specifically, our CGPN derives from variational inference; integrates a newly designed Graph Poisson Network (GPN) to effectively propagate the limited labels to the entire graph and a normal GNN, such as Graph Attention Network, that flexibly guides the propagation of GPN; applies a contrastive objective to further exploit the supervision information from the learning process of GPN and GNN models.

Essentially, our CGPN can enhance the learning performance of GNNs under extremely limited labels by contrastively propagating the limited labels to the entire graph. We conducted extensive experiments on different types of datasets to demonstrate the superiority of CGPN.

[Collaborative Uncertainty in Multi-Agent Trajectory Forecasting](#)

- Bohan Tang, Yiqi Zhong, Ulrich Neumann, Gang Wang, Siheng Chen, Ya Zhang
- abstract@[open-review](#): Uncertainty modeling is critical in trajectory-forecasting systems for both interpretation and safety reasons. To better predict the future trajectories of multiple agents, recent works have introduced interaction modules to capture interactions among agents. This approach leads to correlations among the predicted trajectories. However, the uncertainty brought by such correlations is neglected. To fill this gap, we propose a novel concept, collaborative uncertainty (CU), which models the uncertainty resulting from the interaction module. We build a general CU-based framework to make a prediction model learn the future trajectory and the corresponding uncertainty. The CU-based framework is integrated as a plugin module to current state-of-the-art (SOTA) systems and deployed in two special cases based on multivariate Gaussian and Laplace distributions. In each case, we conduct extensive experiments on two synthetic datasets and two public, large-scale benchmarks of trajectory forecasting. The results are promising: 1) The results of synthetic datasets show that CU-based framework allows the model to nicely rebuild the ground-truth distribution. 2) The results of trajectory forecasting benchmarks demonstrate that the CU-based framework steadily helps SOTA systems improve their performances. Specially, the proposed CU-based framework helps VectorNet improve by 57 cm regarding Final Displacement Error on nuScenes dataset. 3) The visualization results of CU illustrate that the value of CU is highly related to the amount of the interactive information among agents.

[Network-to-Network Regularization: Enforcing Occam's Razor to Improve Generalization](#)

- Rohan Ghosh, Mehul Motani
- abstract@[open-review](#): What makes a classifier have the ability to generalize? There have been a lot of important attempts to address this question, but a clear answer is still elusive. Proponents of complexity theory find that the complexity of the classifier's function space is key to deciding generalization, whereas other recent work reveals that classifiers which extract invariant feature representations are likely to generalize better. Recent theoretical and empirical studies, however, have shown that even within a classifier's function space, there can be significant differences in the ability to generalize. Specifically, empirical studies have shown that among functions which have a good training data fit, functions with lower Kolmogorov complexity (KC) are likely to generalize better, while the opposite is true for functions of higher KC. Motivated by these findings, we propose, in this work, a novel measure of complexity called Kolmogorov Growth (KG), which we use to derive new generalization error bounds that only depend on the final choice of the classification function. Guided by the bounds, we propose a novel way of regularizing neural networks by constraining the network trajectory to remain in the low KG zone during training. Minimizing KG while learning is akin to applying the Occam's razor to neural networks. The proposed approach, called network-to-network regularization, leads to clear improvements in the generalization ability of classifiers. We verify this for three popular image datasets (MNIST, CIFAR-10, CIFAR-100) across varying training data sizes. Empirical studies find that conventional training of neural networks, unlike network-to-network regularization, leads to networks of high KG and lower test accuracies. Furthermore, we present the benefits of N2N regularization in the scenario where the training data labels are noisy. Using N2N regularization, we achieve competitive performance on MNIST, CIFAR-10 and CIFAR-100 datasets with corrupted training labels, significantly improving network performance compared to standard cross-entropy baselines in most cases. These findings illustrate the many benefits obtained from imposing a function complexity prior like Kolmogorov Growth during the training process.

[Generalized and Discriminative Few-Shot Object Detection via SVD-Dictionary Enhancement](#)

- Aming WU, Suqi Zhao, Cheng Deng, Wei Liu
- abstract@[open-review](#): Few-shot object detection (FSOD) aims to detect new objects based on few annotated samples. To alleviate the impact of few samples, enhancing the generalization and discrimination abilities of detectors on new objects plays an important role. In this paper, we explore employing Singular Value Decomposition (SVD) to boost both the generalization and discrimination abilities. In specific, we propose a novel method, namely, SVD-Dictionary enhancement, to build two separated spaces based on the sorted singular values. Concretely, the eigenvectors corresponding to larger singular values are used to build the generalization space in which localization is performed, as these eigenvectors generally suppress certain variations (e.g., the variation of styles) and contain intrinsical characteristics of objects. Meanwhile, since the eigenvectors corresponding to relatively smaller singular values may contain richer category-related information, we can utilize them to build the discrimination space in which classification is performed. Dictionary learning is further leveraged to capture high-level discriminative information from the discrimination space, which is beneficial for improving detection accuracy. In the experiments, we separately verify the effectiveness of our method on PASCAL VOC and COCO benchmarks. Particularly, for the 2-shot case in VOC split1, our method significantly outperforms the baseline by 6.2%. Moreover, visualization analysis shows that our method is instrumental in doing FSOD.

[Conditioning Sparse Variational Gaussian Processes for Online Decision-making](#)

- Wesley J. Maddox, Samuel Stanton, Andrew G. Wilson
- abstract@[open-review](#): With a principled representation of uncertainty and closed form posterior updates, Gaussian processes (GPs) are a natural choice for online decision making. However, Gaussian processes typically require at least $\mathcal{O}(n^2)$ computations for n training points, limiting their general applicability. Stochastic variational Gaussian processes (SVGPs) can provide scalable inference for a dataset of fixed size, but are difficult to efficiently condition on new data. We propose online variational conditioning (OVC), a procedure for efficiently conditioning SVGPs in an online setting that does not require re-training through the evidence lower bound with the addition of new data. OVC enables the pairing of SVGPs with advanced look-ahead acquisition functions for black-box optimization, even with non-Gaussian likelihoods. We show OVC provides compelling performance in a range of applications including active learning of malaria incidence, and reinforcement learning on MuJoCo simulated robotic control tasks.

[Spherical Motion Dynamics: Learning Dynamics of Normalized Neural Network using SGD and Weight Decay](#)

- Ruosi Wan, Zhanxing Zhu, Xiangyu Zhang, Jian Sun
- abstract@[open-review](#): In this paper, we comprehensively reveal the learning dynamics of normalized neural network using Stochastic Gradient Descent (with momentum) and Weight Decay (WD), named as Spherical Motion Dynamics (SMD). Most related works focus on studying behavior of effective learning rate "inequilibrium" state, i.e. assuming weight norm remains unchanged. However, their discussion on why this equilibrium can be reached is either absent or less convincing. Our work directly explores the cause of equilibrium, as a special state of SMD. Specifically, 1) we introduce the assumptions that can lead to equilibrium state in SMD, and prove equilibrium can be reached in a linear rate regime under given assumptions; 2) we propose ``angular update'' as a substitute for effective learning rate to depict the state of SMD, and derive the theoretical value of angular update in equilibrium state; 3) we verify our assumptions and theoretical results on various large-scale computer vision tasks including ImageNet and MSCOCO with standard settings. Experiment results show our theoretical findings agree well with empirical observations. We also show that the behavior of angular update in SMD can produce interesting effect to the optimization of neural network in practice.

[Imitating Deep Learning Dynamics via Locally Elastic Stochastic Differential Equations](#)

- Jiayao Zhang, Hua Wang, Weijie Su
- abstract@[open-review](#): Understanding the training dynamics of deep learning models is perhaps a necessary step toward demystifying the effectiveness of these models. In particular, how do training data from different classes gradually become separable in their feature spaces when training neural networks

using stochastic gradient descent? In this paper, we model the evolution of features during deep learning training using a set of stochastic differential equations (SDEs) that each corresponding to a training sample. As a crucial ingredient in our modeling strategy, each SDE contains a drift term that reflects the impact of backpropagation at an input on the features of all samples. Our main finding uncovers a sharp phase transition phenomenon regarding the intra-class impact: if the SDEs are locally elastic in the sense that the impact is more significant on samples from the same class as the input, the features of training data become linearly separable---meaning vanishing training loss; otherwise, the features are not separable, no matter how long the training time is. In the presence of local elasticity, moreover, an analysis of our SDEs shows the emergence of a simple geometric structure called neural collapse of the features. Taken together, our results shed light on the decisive role of local elasticity underlying the training dynamics of neural networks. We corroborate our theoretical analysis with experiments on a synthesized dataset of geometric shapes as well as on CIFAR-10.

[Probabilistic Forecasting: A Level-Set Approach](#)

- Hilaf Hasson, Bernie Wang, Tim Januschowski, Jan Gasthaus
- abstract@[open-review](#): Large-scale time series panels have become ubiquitous over the last years in areas such as retail, operational metrics, IoT, and medical domain (to name only a few). This has resulted in a need for forecasting techniques that effectively leverage all available data by learning across all time series in each panel. Among the desirable properties of forecasting techniques, being able to generate probabilistic predictions ranks among the top. In this paper, we therefore present Level Set Forecaster (LSF), a simple yet effective general approach to transform a point estimator into a probabilistic one. By recognizing the connection of our algorithm to random forests (RFs) and quantile regression forests (QRFs), we are able to prove consistency guarantees of our approach under mild assumptions on the underlying point estimator. As a byproduct, we prove the first consistency results for QRFs under the CART-splitting criterion. Empirical experiments show that our approach, equipped with tree-based models as the point estimator, rivals state-of-the-art deep learning models in terms of forecasting accuracy.

[Roto-translated Local Coordinate Frames For Interacting Dynamical Systems](#)

- Miltiadis Kofinas, Naveen Nagaraja, Efstratios Gavves
- abstract@[open-review](#): Modelling interactions is critical in learning complex dynamical systems, namely systems of interacting objects with highly non-linear and time-dependent behaviour. A large class of such systems can be formalized as $\{\text{geometric graphs}\}$, $\{\text{i.e.}\}$ graphs with nodes positioned in the Euclidean space given an $\{\text{arbitrarily}\}$ chosen global coordinate system, for instance vehicles in a traffic scene. Notwithstanding the arbitrary global coordinate system, the governing dynamics of the respective dynamical systems are invariant to rotations and translations, also known as $\{\text{Galilean invariance}\}$. As ignoring these invariances leads to worse generalization, in this work we propose local coordinate systems per node-object to induce roto-translation invariance to the geometric graph of the interacting dynamical system. Further, the local coordinate systems allow for a natural definition of anisotropic filtering in graph neural networks. Experiments in traffic scenes, 3D motion capture, and colliding particles demonstrate the proposed approach comfortably outperforms the recent state-of-the-art.

[ParK: Sound and Efficient Kernel Ridge Regression by Feature Space Partitions](#)

- Luigi Carratino, Stefano Vigogna, Daniele Calandriello, Lorenzo Rosasco
- abstract@[open-review](#): We introduce ParK, a new large-scale solver for kernel ridge regression. Our approach combines partitioning with random projections and iterative optimization to reduce space and time complexity while provably maintaining the same statistical accuracy. In particular, constructing suitable partitions directly in the feature space rather than in the input space, we promote orthogonality between the local estimators, thus ensuring that key quantities such as local effective dimension and bias remain under control. We characterize the statistical-computational tradeoff of our model, and demonstrate the effectiveness of our method by numerical experiments on large-scale datasets.

[Scaling Gaussian Processes with Derivative Information Using Variational Inference](#)

- Misha Padidar, Xinran Zhu, Leo Huang, Jacob Gardner, David Bindel
- abstract@[open-review](#): Gaussian processes with derivative information are useful in many settings where derivative information is available, including numerous Bayesian optimization and regression tasks that arise in the natural sciences. Incorporating derivative observations, however, comes with a dominating $O(N^3D^3)$ computational cost when training on N points in D input dimensions. This is intractable for even moderately sized problems. While recent work has addressed this intractability in the low- D setting, the high- N , high- D setting is still unexplored and of great value, particularly as machine learning problems increasingly become high dimensional. In this paper, we introduce methods to achieve fully scalable Gaussian process regression with derivatives using variational inference. Analogous to the use of inducing values to sparsify the labels of a training set, we introduce the concept of inducing directional derivatives to sparsify the partial derivative information of the training set. This enables us to construct a variational posterior that incorporates derivative information but whose size depends neither on the full dataset size N nor the full dimensionality D . We demonstrate the full scalability of our approach on a variety of tasks, ranging from a high dimensional Stellarator fusion regression task to training graph convolutional neural networks on PubMed using Bayesian optimization. Surprisingly, we additionally find that our approach can improve regression performance even in settings where only label data is available.

[On the Representation of Solutions to Elliptic PDEs in Barron Spaces](#)

- Ziang Chen, Jianfeng Lu, Yulong Lu
- abstract@[open-review](#): Numerical solutions to high-dimensional partial differential equations (PDEs) based on neural networks have seen exciting developments. This paper derives complexity estimates of the solutions of d -dimensional second-order elliptic PDEs in the Barron space, that is a set of functions admitting the integral of certain parametric ridge function against a probability measure on the parameters. We prove under some appropriate assumptions that if the coefficients and the source term of the elliptic PDE lie in Barron spaces, then the solution of the PDE is ϵ -close with respect to the H^1 norm to a Barron function. Moreover, we prove dimension-explicit bounds for the Barron norm of this approximate solution, depending at most polynomially on the dimension d of the PDE. As a direct consequence of the complexity estimates, the solution of the PDE can be approximated on any bounded domain by a two-layer neural network with respect to the H^1 norm with a dimension-explicit convergence rate.

[A/B Testing for Recommender Systems in a Two-sided Marketplace](#)

- Preetam Nandy, Divya Venugopalan, Chun Lo, Shaunak Chatterjee
- abstract@[open-review](#): Two-sided marketplaces are standard business models of many online platforms (e.g., Amazon, Facebook, LinkedIn), wherein the platforms have consumers, buyers or content viewers on one side and producers, sellers or content-creators on the other. Consumer side measurement of the impact of a treatment variant can be done via simple online A/B testing. Producer side measurement is more challenging because the producer experience depends on the treatment assignment of the consumers. Existing approaches for producer side measurement are either based on graph cluster-based randomization or on certain treatment propagation assumptions. The former approach results in low-powered experiments as the producer-consumer network density increases and the latter approach lacks a strict notion of error control. In this paper, we propose (i) a quantification of the quality of a producer side experiment design, and (ii) a new experiment design mechanism that generates high-quality experiments based on this quantification. Our approach, called UniCoRn (Unifying Counterfactual Rankings), provides explicit control over the quality of the experiment and its computation cost. Further, we prove that our experiment design is optimal to the proposed design quality measure. Our approach is agnostic to the density of the producer-consumer network and does not rely on any treatment propagation assumption. Moreover, unlike the existing approaches, we do not need to know the underlying network in advance, making this widely applicable to the industrial setting where the underlying network is unknown and challenging to

predict a priori due to its dynamic nature. We use simulations to validate our approach and compare it against existing methods. We also deployed UniCoRN in an edge recommendation application that serves tens of millions of members and billions of edge recommendations daily.

[Retiring Adult: New Datasets for Fair Machine Learning](#)

- Frances Ding, Moritz Hardt, John Miller, Ludwig Schmidt
- abstract@[open-review](#): Although the fairness community has recognized the importance of data, researchers in the area primarily rely on UCI Adult when it comes to tabular data. Derived from a 1994 US Census survey, this dataset has appeared in hundreds of research papers where it served as the basis for the development and comparison of many algorithmic fairness interventions. We reconstruct a superset of the UCI Adult data from available US Census sources and reveal idiosyncrasies of the UCI Adult dataset that limit its external validity. Our primary contribution is a suite of new datasets derived from US Census surveys that extend the existing data ecosystem for research on fair machine learning. We create prediction tasks relating to income, employment, health, transportation, and housing. The data span multiple years and all states of the United States, allowing researchers to study temporal shift and geographic variation. We highlight a broad initial sweep of new empirical insights relating to trade-offs between fairness criteria, performance of algorithmic interventions, and the role of distribution shift based on our new datasets. Our findings inform ongoing debates, challenge some existing narratives, and point to future research directions.

[Cardinality constrained submodular maximization for random streams](#)

- Paul Liu, Aviad Rubinstein, Jan Vondrak, Junyao Zhao
- abstract@[open-review](#): We consider the problem of maximizing submodular functions in single-pass streaming and secretaries-with-shortlists models, both with random arrival order. For cardinality constrained monotone functions, Agrawal, Shadravan, and Stein~\cite{SMC19} gave a single-pass $(1-1/e-\epsilon)$ -approximation algorithm using only linear memory, but their exponential dependence on ϵ makes it impractical even for $\epsilon=0.1$. We simplify both the algorithm and the analysis, obtaining an exponential improvement in the ϵ -dependence (in particular, $O(k/\epsilon)$ memory). Extending these techniques, we also give a simple $(1/e-\epsilon)$ -approximation for non-monotone functions in $O(k/\epsilon)$ memory. For the monotone case, we also give a corresponding unconditional hardness barrier of $1-1/e+\epsilon$ for single-pass algorithms in randomly ordered streams, even assuming unlimited computation. Finally, we show that the algorithms are simple to implement and work well on real world datasets.

[Self-Instantiated Recurrent Units with Dynamic Soft Recursion](#)

- Aston Zhang, Yi Tay, Yikang Shen, Alvin Chan, SHUAI ZHANG
- abstract@[open-review](#): While standard recurrent neural networks explicitly impose a chain structure on different forms of data, they do not have an explicit bias towards recursive self-instantiation where the extent of recursion is dynamic. Given diverse and even growing data modalities (e.g., logic, algorithmic input and output, music, code, images, and language) that can be expressed in sequences and may benefit from more architectural flexibility, we propose the self-instantiated recurrent unit (Self-IRU) with a novel inductive bias towards dynamic soft recursion. On one hand, the Self-IRU is characterized by recursive self-instantiation via its gating functions, i.e., gating mechanisms of the Self-IRU are controlled by instances of the Self-IRU itself, which are repeatedly invoked in a recursive fashion. On the other hand, the extent of the Self-IRU recursion is controlled by gates whose values are between 0 and 1 and may vary across the temporal dimension of sequences, enabling dynamic soft recursion depth at each time step. The architectural flexibility and effectiveness of our proposed approach are demonstrated across multiple data modalities. For example, the Self-IRU achieves state-of-the-art performance on the logical inference dataset [Bowman et al., 2014] even when comparing with competitive models that have access to ground-truth syntactic information.

[Sparse Uncertainty Representation in Deep Learning with Inducing Weights](#)

- Hippolyt Ritter, Martin Kukla, Cheng Zhang, Yingzhen Li
- abstract@[open-review](#): Bayesian Neural Networks and deep ensembles represent two modern paradigms of uncertainty quantification in deep learning. Yet these approaches struggle to scale mainly due to memory inefficiency, requiring parameter storage several times that of their deterministic counterparts. To address this, we augment each weight matrix with a small inducing weight matrix, projecting the uncertainty quantification into a lower dimensional space. We further extend Matheron's conditional Gaussian sampling rule to enable fast weight sampling, which enables our inference method to maintain reasonable run-time as compared with ensembles. Importantly, our approach achieves competitive performance to the state-of-the-art in prediction and uncertainty estimation tasks with fully connected neural networks and ResNets, while reducing the parameter size to $\leq 24.3\%$ of that of a single neural network.

[Scalable Inference of Sparsely-changing Gaussian Markov Random Fields](#)

- Salar Fattah, Andres Gomez
- abstract@[open-review](#): We study the problem of inferring time-varying Gaussian Markov random fields, where the underlying graphical model is both sparse and changes {sparsely} over time. Most of the existing methods for the inference of time-varying Markov random fields (MRFs) rely on the regularized maximum likelihood estimation (MLE), that typically suffer from weak statistical guarantees and high computational time. Instead, we introduce a new class of constrained optimization problems for the inference of sparsely-changing Gaussian MRFs (GMRFs). The proposed optimization problem is formulated based on the exact ℓ_0 regularization, and can be solved in near-linear time and memory. Moreover, we show that the proposed estimator enjoys a provably small estimation error. We derive sharp statistical guarantees in the high-dimensional regime, showing that such problems can be learned with as few as one sample per time period. Our proposed method is extremely efficient in practice: it can accurately estimate sparsely-changing GMRFs with more than 500 million variables in less than one hour.

[Grad2Task: Improved Few-shot Text Classification Using Gradients for Task Representation](#)

- Jixuan Wang, Kuan-Chieh Wang, Frank Rudzicz, Michael Brudno
- abstract@[open-review](#): Large pretrained language models (LMs) like BERT have improved performance in many disparate natural language processing (NLP) tasks. However, fine tuning such models requires a large number of training examples for each target task. Simultaneously, many realistic NLP problems are "few shot", without a sufficiently large training set. In this work, we propose a novel conditional neural process-based approach for few-shot text classification that learns to transfer from other diverse tasks with rich annotation. Our key idea is to represent each task using gradient information from a base model and to train an adaptation network that modulates a text classifier conditioned on the task representation. While previous task-aware few-shot learners represent tasks by input encoding, our novel task representation is more powerful, as the gradient captures input-output relationships of a task. Experimental results show that our approach outperforms traditional fine-tuning, sequential transfer learning, and state-of-the-art meta learning approaches on a collection of diverse few-shot tasks. We further conducted analysis and ablations to justify our design choices.

[Learnability of Linear Thresholds from Label Proportions](#)

- Rishi Saket

- abstract@[open-review](#): We study the problem of properly learning linear threshold functions (LTFs) in the learning from label proportions (LLP) framework. In this, the learning is on a collection of bags of feature-vectors with only the proportion of labels available for each bag. First, we provide an algorithm that, given a collection of such bags each of size at most two whose label proportions are consistent with (i.e., the bags are satisfied by) an unknown LTF, efficiently produces an LTF that satisfies at least \$(2/5)\$-fraction of the bags. If all the bags are non-monochromatic (i.e., bags of size two with differently labeled feature-vectors) the algorithm satisfies at least \$(1/2)\$-fraction of them. For the special case of OR over the \$d\$-dimensional boolean vectors, we give an algorithm which computes an LTF achieving an additional \$\Omega(1/d)\$ in accuracy for the two cases. Our main result provides evidence that these algorithmic bounds cannot be significantly improved, even for learning monotone ORs using LTFs. We prove that it is NP-hard, given a collection of non-monochromatic bags which are all satisfied by some monotone OR, to compute any function of constantly many LTFs that satisfies \$(1/2 + \epsilon)\$-fraction of the bags for any constant \$\epsilon > 0\$. This bound is tight for the non-monochromatic bags case. The above is in contrast to the usual supervised learning setup (i.e., unit-sized bags) in which LTFs are efficiently learnable to arbitrary accuracy using linear programming, and even a trivial algorithm (any LTF or its complement) achieves an accuracy of \$1/2\$. These techniques however, fail in the LLP setting. Indeed, we show that the LLP learning of LTFs (even for the special case of monotone ORs) using LTFs dramatically increases in complexity as soon as bags of size two are allowed. Our work gives the first inapproximability for LLP learning LTFs, and a strong complexity separation between LLP and traditional supervised learning.

[A variational approximate posterior for the deep Wishart process](#)

- Sebastian Ober, Laurence Aitchison
- abstract@[open-review](#): Recent work introduced deep kernel processes as an entirely kernel-based alternative to NNs (Aitchison et al. 2020). Deep kernel processes flexibly learn good top-layer representations by alternately sampling the kernel from a distribution over positive semi-definite matrices and performing nonlinear transformations. A particular deep kernel process, the deep Wishart process (DWP), is of particular interest because its prior can be made equivalent to deep Gaussian process (DGP) priors for kernels that can be expressed entirely in terms of Gram matrices. However, inference in DWPs has not yet been possible due to the lack of sufficiently flexible distributions over positive semi-definite matrices. Here, we give a novel approach to obtaining flexible distributions over positive semi-definite matrices by generalising the Bartlett decomposition of the Wishart probability density. We use this new distribution to develop an approximate posterior for the DWP that includes dependency across layers. We develop a doubly-stochastic inducing-point inference scheme for the DWP and show experimentally that inference in the DWP can improve performance over doing inference in a DGP with the equivalent prior.

[Neural Pseudo-Label Optimism for the Bank Loan Problem](#)

- Aldo Pacchiano, Shaun Singh, Edward Chou, Alex Berg, Jakob Foerster
- abstract@[open-review](#): We study a class of classification problems best exemplified by the \emph{bank loan} problem, where a lender decides whether or not to issue a loan. The lender only observes whether a customer will repay a loan if the loan is issued to begin with, and thus modeled decisions affect what data is available to the lender for future decisions. As a result, it is possible for the lender's algorithm to ``get stuck'' with a self-fulfilling model. This model never corrects its false negatives, since it never sees the true label for rejected data, thus accumulating infinite regret. In the case of linear models, this issue can be addressed by adding optimism directly into the model predictions. However, there are few methods that extend to the function approximation case using Deep Neural Networks. We present Pseudo-Label Optimism (PLOT), a conceptually and computationally simple method for this setting applicable to DNNs. \texttt{PLOT} adds an optimistic label to the subset of decision points the current model is deciding on, trains the model on all data so far (including these points along with their optimistic labels), and finally uses the resulting \emph{optimistic} model for decision making. \texttt{PLOT} achieves competitive performance on a set of three challenging benchmark problems, requiring minimal hyperparameter tuning. We also show that \texttt{PLOT} satisfies a logarithmic regret guarantee, under a Lipschitz and logistic mean label model, and under a separability condition on the data.

[Visualizing the Emergence of Intermediate Visual Patterns in DNNs](#)

- Mingjie Li, Shaobo Wang, Quanshi Zhang
- abstract@[open-review](#): This paper proposes a method to visualize the discrimination power of intermediate-layer visual patterns encoded by a DNN. Specifically, we visualize (1) how the DNN gradually learns regional visual patterns in each intermediate layer during the training process, and (2) the effects of the DNN using non-discriminative patterns in low layers to construct discriminative patterns in middle/high layers through the forward propagation. Based on our visualization method, we can quantify knowledge points (i.e. the number of discriminative visual patterns) learned by the DNN to evaluate the representation capacity of the DNN. Furthermore, this method also provides new insights into signal-processing behaviors of existing deep-learning techniques, such as adversarial attacks and knowledge distillation.

[Learning 3D Dense Correspondence via Canonical Point Autoencoder](#)

- An-Chieh Cheng, Xuetong Li, Min Sun, Ming-Hsuan Yang, Sifei Liu
- abstract@[open-review](#): We propose a canonical point autoencoder (CPAE) that predicts dense correspondences between 3D shapes of the same category. The autoencoder performs two key functions: (a) encoding an arbitrarily ordered point cloud to a canonical primitive, e.g., a sphere, and (b) decoding the primitive back to the original input instance shape. As being placed in the bottleneck, this primitive plays a key role to map all the unordered point clouds on the canonical surface, and to be reconstructed in an ordered fashion. Once trained, points from different shape instances that are mapped to the same locations on the primitive surface are determined to be a pair of correspondence. Our method does not require any form of annotation or self-supervised part segmentation network and can handle unaligned input point clouds within a certain rotation range. Experimental results on 3D semantic keypoint transfer and part segmentation transfer show that our model performs favorably against state-of-the-art correspondence learning methods.

[Speech-T: Transducer for Text to Speech and Beyond](#)

- Jiawei Chen, Xu Tan, Yichong Leng, Jin Xu, Guihua Wen, Tao Qin, Tie-Yan Liu
- abstract@[open-review](#): Neural Transducer (e.g., RNN-T) has been widely used in automatic speech recognition (ASR) due to its capabilities of efficiently modeling monotonic alignments between input and output sequences and naturally supporting streaming inputs. Considering that monotonic alignments are also critical to text to speech (TTS) synthesis and streaming TTS is also an important application scenario, in this work, we explore the possibility of applying Transducer to TTS and more. However, it is challenging because it is difficult to trade off the emission (continuous mel-spectrogram prediction) probability and transition (ASR Transducer predicts blank token to indicate transition to next input) probability when calculating the output probability lattice in Transducer, and it is not easy to learn the alignments between text and speech through the output probability lattice. We propose SpeechTransducer (Speech-T for short), a Transformer based Transducer model that 1) uses a new forward algorithm to separate the transition prediction from the continuous mel-spectrogram prediction when calculating the output probability lattice, and uses a diagonal constraint in the probability lattice to help the alignment learning; 2) supports both full-sentence or streaming TTS by adjusting the look-ahead context; and 3) further supports both TTS and ASR together for the first time, which enjoys several advantages including fewer parameters as well as streaming synthesis and recognition in a single model. Experiments on LJSpeech datasets demonstrate that Speech-T 1) is more robust than the attention based autoregressive TTS model due to its inherent monotonic alignments between text and speech; 2) naturally supports streaming TTS with good voice quality; and 3) enjoys the benefit of joint modeling TTS and ASR in a single network.

[Multi-modal Dependency Tree for Video Captioning](#)

- Wentian Zhao, Xinxiao Wu, Jiebo Luo
- abstract@[open-review](#): Generating fluent and relevant language to describe visual content is critical for the video captioning task. Many existing methods generate captions using sequence models that predict words in a left-to-right order. In this paper, we investigate a graph-structured model for caption generation by explicitly modeling the hierarchical structure in the sentences to further improve the fluency and relevance of sentences. To this end, we propose a novel video captioning method that generates a sentence by first constructing a multi-modal dependency tree and then traversing the constructed tree, where the syntactic structure and semantic relationship in the sentence are represented by the tree topology. To take full advantage of the information from both vision and language, both the visual and textual representation features are encoded into each tree node. Different from existing dependency parsing methods that generate uni-modal dependency trees for language understanding, our method constructs multi-modal dependency trees for language generation of images and videos. We also propose a tree-structured reinforcement learning algorithm to effectively optimize the captioning model where a novel reward is designed by evaluating the semantic consistency between the generated sub-tree and the ground-truth tree. Extensive experiments on several video captioning datasets demonstrate the effectiveness of the proposed method.

[Greedy and Random Quasi-Newton Methods with Faster Explicit Superlinear Convergence](#)

- Dachao Lin, Haishan Ye, Zhihua Zhang
- abstract@[open-review](#): In this paper, we follow Rodomanov and Nesterov's work to study quasi-Newton methods. We focus on the common SR1 and BFGS quasi-Newton methods to establish better explicit (local) superlinear convergence rates. First, based on the greedy quasi-Newton update which greedily selects the direction to maximize a certain measure of progress, we improve the convergence rate to a condition-number-free superlinear convergence rate. Second, based on the random quasi-Newton update that selects the direction randomly from a spherically symmetric distribution, we show the same superlinear convergence rate established as above. Our analysis is closely related to the approximation of a given Hessian matrix, unconstrained quadratic objective, as well as the general strongly convex, smooth, and strongly self-concordant functions.

[Neural Tangent Kernel Maximum Mean Discrepancy](#)

- Xiuyuan Cheng, Yao Xie
- abstract@[open-review](#): We present a novel neural network Maximum Mean Discrepancy (MMD) statistic by identifying a new connection between neural tangent kernel (NTK) and MMD. This connection enables us to develop a computationally efficient and memory-efficient approach to compute the MMD statistic and perform NTK based two-sample tests towards addressing the long-standing challenge of memory and computational complexity of the MMD statistic, which is essential for online implementation to assimilating new samples. Theoretically, such a connection allows us to understand the NTK test statistic properties, such as the Type-I error and testing power for performing the two-sample test, by adapting existing theories for kernel MMD. Numerical experiments on synthetic and real-world datasets validate the theory and demonstrate the effectiveness of the proposed NTK-MMD statistic.

[Subgraph Federated Learning with Missing Neighbor Generation](#)

- Ke ZHANG, Carl Yang, Xiaoxiao Li, Lichao Sun, Siu Ming Yiu
- abstract@[open-review](#): Graphs have been widely used in data mining and machine learning due to their unique representation of real-world objects and their interactions. As graphs are getting bigger and bigger nowadays, it is common to see their subgraphs separately collected and stored in multiple local systems. Therefore, it is natural to consider the subgraph federated learning setting, where each local system holds a small subgraph that may be biased from the distribution of the whole graph. Hence, the subgraph federated learning aims to collaboratively train a powerful and generalizable graph mining model without directly sharing their graph data. In this work, towards the novel yet realistic setting of subgraph federated learning, we propose two major techniques: (1) FedSage, which trains a GraphSage model based on FedAvg to integrate node features, link structures, and task labels on multiple local subgraphs; (2) FedSage+, which trains a missing neighbor generator along FedSage to deal with missing links across local subgraphs. Empirical results on four real-world graph datasets with synthesized subgraph federated learning settings demonstrate the effectiveness and efficiency of our proposed techniques. At the same time, consistent theoretical implications are made towards their generalization ability on the global graphs.

[Bellman-consistent Pessimism for Offline Reinforcement Learning](#)

- Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, Alekh Agarwal
- abstract@[open-review](#): The use of pessimism, when reasoning about datasets lacking exhaustive exploration has recently gained prominence in offline reinforcement learning. Despite the robustness it adds to the algorithm, overly pessimistic reasoning can be equally damaging in precluding the discovery of good policies, which is an issue for the popular bonus-based pessimism. In this paper, we introduce the notion of Bellman-consistent pessimism for general function approximation: instead of calculating a point-wise lower bound for the value function, we implement pessimism at the initial state over the set of functions consistent with the Bellman equations. Our theoretical guarantees only require Bellman closedness as standard in the exploratory setting, in which case bonus-based pessimism fails to provide guarantees. Even in the special case of linear function approximation where stronger expressivity assumptions hold, our result improves upon a recent bonus-based approach by $\mathcal{O}(d)$ in its sample complexity (when the action space is finite). Remarkably, our algorithms automatically adapt to the best bias-variance tradeoff in the hindsight, whereas most prior approaches require tuning extra hyperparameters a priori.

[Can You Learn an Algorithm? Generalizing from Easy to Hard Problems with Recurrent Networks](#)

- Avi Schwarzschild, Eitan Borgnia, Arjun Gupta, Furong Huang, Uzi Vishkin, Micah Goldblum, Tom Goldstein
- abstract@[open-review](#): Deep neural networks are powerful machines for visual pattern recognition, but reasoning tasks that are easy for humans may still be difficult for neural models. Humans possess the ability to extrapolate reasoning strategies learned on simple problems to solve harder examples, often by thinking for longer. For example, a person who has learned to solve small mazes can easily extend the very same search techniques to solve much larger mazes by spending more time. In computers, this behavior is often achieved through the use of algorithms, which scale to arbitrarily hard problem instances at the cost of more computation. In contrast, the sequential computing budget of feed-forward neural networks is limited by their depth, and networks trained on simple problems have no way of extending their reasoning to accommodate harder problems. In this work, we show that recurrent networks trained to solve simple problems with few recurrent steps can indeed solve much more complex problems simply by performing additional recurrences during inference. We demonstrate this algorithmic behavior of recurrent networks on prefix sum computation, mazes, and chess. In all three domains, networks trained on simple problem instances are able to extend their reasoning abilities at test time simply by "thinking for longer."

[Sub-Linear Memory: How to Make Performers SLiM](#)

- Valerii Likhoshesterov, Krzysztof M. Choromanski, Jared Quincy Davis, Xingyou Song, Adrian Weller
- abstract@[open-review](#): Transformer architectures have become very popular yet the original implementation requires $\mathcal{O}(L^2)$ in serial time and memory as functions of input length L . Recent works proposed various linear self-attention mechanisms, scaling only as $\mathcal{O}(L)$ for serial computation. We conduct a thorough complexity analysis of Performers, a class which includes most recent linear Transformer mechanisms. We note a remarkable computational flexibility: the gradient computation can be performed with no approximations using sublinear memory as a function of L (in addition to negligible storage for the input sequence), at a cost of greater time complexity in the parallel setting. In the extreme case, a Performer consumes only $\mathcal{O}(1)$ memory, and still requires $\mathcal{O}(L)$ time. Due to complete backward-compatibility, this discovered time-memory tradeoff can be used for fine-tuning on low-memory devices in a decentralized fashion without any server computations.

Efficient Learning of Discrete-Continuous Computation Graphs

- David Friede, Mathias Niepert
- abstract@[open-review](#): Numerous models for supervised and reinforcement learning benefit from combinations of discrete and continuous model components. End-to-end learnable discrete-continuous models are compositional, tend to generalize better, and are more interpretable. A popular approach to building discrete-continuous computation graphs is that of integrating discrete probability distributions into neural networks using stochastic softmax tricks. Prior work has mainly focused on computation graphs with a single discrete component on each of the graph's execution paths. We analyze the behavior of more complex stochastic computations graphs with multiple sequential discrete components. We show that it is challenging to optimize the parameters of these models, mainly due to small gradients and local minima. We then propose two new strategies to overcome these challenges. First, we show that increasing the scale parameter of the Gumbel noise perturbations during training improves the learning behavior. Second, we propose dropout residual connections specifically tailored to stochastic, discrete-continuous computation graphs. With an extensive set of experiments, we show that we can train complex discrete-continuous models which one cannot train with standard stochastic softmax tricks. We also show that complex discrete-stochastic models generalize better than their continuous counterparts on several benchmark datasets.

VQ-GNN: A Universal Framework to Scale up Graph Neural Networks using Vector Quantization

- Mucong Ding, Kezhi Kong, Jingling Li, Chen Zhu, John Dickerson, Furong Huang, Tom Goldstein
- abstract@[open-review](#): Most state-of-the-art Graph Neural Networks (GNNs) can be defined as a form of graph convolution which can be realized by message passing between direct neighbors or beyond. To scale such GNNs to large graphs, various neighbor-, layer-, or subgraph-sampling techniques are proposed to alleviate the "neighbor explosion" problem by considering only a small subset of messages passed to the nodes in a mini-batch. However, sampling-based methods are difficult to apply to GNNs that utilize many-hops-away or global context each layer, show unstable performance for different tasks and datasets, and do not speed up model inference. We propose a principled and fundamentally different approach, VQ-GNN, a universal framework to scale up any convolution-based GNNs using Vector Quantization (VQ) without compromising the performance. In contrast to sampling-based techniques, our approach can effectively preserve all the messages passed to a mini-batch of nodes by learning and updating a small number of quantized reference vectors of global node representations, using VQ within each GNN layer. Our framework avoids the "neighbor explosion" problem of GNNs using quantized representations combined with a low-rank version of the graph convolution matrix. We show that such a compact low-rank version of the gigantic convolution matrix is sufficient both theoretically and experimentally. In company with VQ, we design a novel approximated message passing algorithm and a nontrivial back-propagation rule for our framework. Experiments on various types of GNN backbones demonstrate the scalability and competitive performance of our framework on large-graph node classification and link prediction benchmarks.

Overcoming Catastrophic Forgetting in Incremental Few-Shot Learning by Finding Flat Minima

- Guangyuan SHI, JIAJIN CHEN, Wenlong Zhang, Li-Ming Zhan, Xiao-Ming Wu
- abstract@[open-review](#): This paper considers incremental few-shot learning, which requires a model to continually recognize new categories with only a few examples provided. Our study shows that existing methods severely suffer from catastrophic forgetting, a well-known problem in incremental learning, which is aggravated due to data scarcity and imbalance in the few-shot setting. Our analysis further suggests that to prevent catastrophic forgetting, actions need to be taken in the primitive stage -- the training of base classes instead of later few-shot learning sessions. Therefore, we propose to search for flat local minima of the base training objective function and then fine-tune the model parameters within the flat region on new tasks. In this way, the model can efficiently learn new classes while preserving the old ones. Comprehensive experimental results demonstrate that our approach outperforms all prior state-of-the-art methods and is very close to the approximate upper bound. The source code is available at <https://github.com/moukamisama/F2M>.

Functional Neural Networks for Parametric Image Restoration Problems

- Fangzhou Luo, Xiaolin Wu, Yanhui Guo
- abstract@[open-review](#): Almost every single image restoration problem has a closely related parameter, such as the scale factor in super-resolution, the noise level in image denoising, and the quality factor in JPEG deblocking. Although recent studies on image restoration problems have achieved great success due to the development of deep neural networks, they handle the parameter involved in an unsophisticated way. Most previous researchers either treat problems with different parameter levels as independent tasks, and train a specific model for each parameter level; or simply ignore the parameter, and train a single model for all parameter levels. The two popular approaches have their own shortcomings. The former is inefficient in computing and the latter is ineffective in performance. In this work, we propose a novel system called functional neural network (FuncNet) to solve a parametric image restoration problem with a single model. Unlike a plain neural network, the smallest conceptual element of our FuncNet is no longer a floating-point variable, but a function of the parameter of the problem. This feature makes it both efficient and effective for a parametric problem. We apply FuncNet to super-resolution, image denoising, and JPEG deblocking. The experimental results show the superiority of our FuncNet on all three parametric image restoration tasks over the state of the arts.

Intrinsic Dimension, Persistent Homology and Generalization in Neural Networks

- Tolga Birdal, Aaron Lou, Leonidas J. Guibas, Umut Simsekli
- abstract@[open-review](#): Disobeying the classical wisdom of statistical learning theory, modern deep neural networks generalize well even though they typically contain millions of parameters. Recently, it has been shown that the trajectories of iterative optimization algorithms can possess fractal structures, and their generalization error can be formally linked to the complexity of such fractals. This complexity is measured by the fractal's intrinsic dimension, a quantity usually much smaller than the number of parameters in the network. Even though this perspective provides an explanation for why overparametrized networks would not overfit, computing the intrinsic dimension (eg, for monitoring generalization during training) is a notoriously difficult task, where existing methods typically fail even in moderate ambient dimensions. In this study, we consider this problem from the lens of topological data analysis (TDA) and develop a generic computational tool that is built on rigorous mathematical foundations. By making a novel connection between learning theory and TDA, we first illustrate that the generalization error can be equivalently bounded in terms of a notion called the 'persistent homology dimension' (PHD), where, compared with prior work, our approach does not require any additional geometrical or statistical assumptions on the training dynamics. Then, by utilizing recently established theoretical results and TDA tools, we develop an efficient algorithm to estimate PHD in the scale of modern deep neural networks and further provide visualization tools to help understand generalization in deep learning. Our experiments show that the proposed approach can efficiently compute a network's intrinsic dimension in a variety of settings, which is predictive of the generalization error.

GemNet: Universal Directional Graph Neural Networks for Molecules

- Johannes Gasteiger, Florian Becker, Stephan Günnemann
- abstract@[open-review](#): Effectively predicting molecular interactions has the potential to accelerate molecular dynamics by multiple orders of magnitude and thus revolutionize chemical simulations. Graph neural networks (GNNs) have recently shown great successes for this task, overtaking classical methods based on fixed molecular kernels. However, they still appear very limited from a theoretical perspective, since regular GNNs cannot distinguish certain types of graphs. In this work we close this gap between theory and practice. We show that GNNs with directed edge embeddings and two-hop message passing are indeed universal approximators for predictions that are invariant to translation, and equivariant to permutation and rotation. We then leverage these insights and multiple structural improvements to propose the geometric message passing neural network (GemNet). We demonstrate the

benefits of the proposed changes in multiple ablation studies. GemNet outperforms previous models on the COLL, MD17, and OC20 datasets by 34%, 41%, and 20%, respectively, and performs especially well on the most challenging molecules. Our implementation is available online.

[Loss function based second-order Jensen inequality and its application to particle variational inference](#)

- Futoshi Futami, Tomoharu Iwata, naonori ueda, Issei Sato, Masashi Sugiyama
- abstract@[open-review](#): Bayesian model averaging, obtained as the expectation of a likelihood function by a posterior distribution, has been widely used for prediction, evaluation of uncertainty, and model selection. Various approaches have been developed to efficiently capture the information in the posterior distribution; one such approach is the optimization of a set of models simultaneously with interaction to ensure the diversity of the individual models in the same way as ensemble learning. A representative approach is particle variational inference (PVI), which uses an ensemble of models as an empirical approximation for the posterior distribution. PVI iteratively updates each model with a repulsion force to ensure the diversity of the optimized models. However, despite its promising performance, a theoretical understanding of this repulsion and its association with the generalization ability remains unclear. In this paper, we tackle this problem in light of PAC-Bayesian analysis. First, we provide a new second-order Jensen inequality, which has the repulsion term based on the loss function. Thanks to the repulsion term, it is tighter than the standard Jensen inequality. Then, we derive a novel generalization error bound and show that it can be reduced by enhancing the diversity of models. Finally, we derive a new PVI that optimizes the generalization error bound directly. Numerical experiments demonstrate that the performance of the proposed PVI compares favorably with existing methods in the experiment.

[Detecting and Adapting to Irregular Distribution Shifts in Bayesian Online Learning](#)

- Aodong Li, Alex Boyd, Padhraic Smyth, Stephan Mandt
- abstract@[open-review](#): We consider the problem of online learning in the presence of distribution shifts that occur at an unknown rate and of unknown intensity. We derive a new Bayesian online inference approach to simultaneously infer these distribution shifts and adapt the model to the detected changes by integrating ideas from change point detection, switching dynamical systems, and Bayesian online learning. Using a binary ‘change variable,’ we construct an informative prior such that--if a change is detected--the model partially erases the information of past model updates by tempering to facilitate adaptation to the new data distribution. Furthermore, the approach uses beam search to track multiple change-point hypotheses and selects the most probable one in hindsight. Our proposed method is model-agnostic, applicable in both supervised and unsupervised learning settings, suitable for an environment of concept drifts or covariate drifts, and yields improvements over state-of-the-art Bayesian online learning approaches.

[Asynchronous Decentralized SGD with Quantized and Local Updates](#)

- Giorgi Nadiradze, Amirmojtaba Sabour, Peter Davies, Shigang Li, Dan Alistarh
- abstract@[open-review](#): Decentralized optimization is emerging as a viable alternative for scalable distributed machine learning, but also introduces new challenges in terms of synchronization costs. To this end, several communication-reduction techniques, such as non-blocking communication, quantization, and local steps, have been explored in the decentralized setting. Due to the complexity of analyzing optimization in such a relaxed setting, this line of work often assumes \emph{global} communication rounds, which require additional synchronization. In this paper, we consider decentralized optimization in the simpler, but harder to analyze, \emph{asynchronous gossip} model, in which communication occurs in discrete, randomly chosen pairings among nodes. Perhaps surprisingly, we show that a variant of SGD called \emph{SwarmSGD} still converges in this setting, even if \emph{non-blocking communication}, \emph{quantization}, and \emph{local steps} are all applied \emph{in conjunction}, and even if the node data distributions and underlying graph topology are both \emph{heterogenous}. Our analysis is based on a new connection with multi-dimensional load-balancing processes. We implement this algorithm and deploy it in a super-computing environment, showing that it can outperform previous decentralized methods in terms of end-to-end training time, and that it can even rival carefully-tuned large-batch SGD for certain tasks.

[Stochastic Shortest Path: Minimax, Parameter-Free and Towards Horizon-Free Regret](#)

- Jean Tarbouriech, Runlong Zhou, Simon S. Du, Matteo Pirotta, Michal Valko, Alessandro Lazaric
- abstract@[open-review](#): We study the problem of learning in the stochastic shortest path (SSP) setting, where an agent seeks to minimize the expected cost accumulated before reaching a goal state. We design a novel model-based algorithm EB-SSP that carefully skews the empirical transitions and perturbs the empirical costs with an exploration bonus to induce an optimistic SSP problem whose associated value iteration scheme is guaranteed to converge. We prove that EB-SSP achieves the minimax regret rate $\$widetilde{O}(B_{\star}\sqrt{SAK})$, where K is the number of episodes, S is the number of states, A is the number of actions and B_{\star} bounds the expected cumulative cost of the optimal policy from any state, thus closing the gap with the lower bound. Interestingly, EB-SSP obtains this result while being parameter-free, i.e., it does not require any prior knowledge of B_{\star} , nor of T_{\star} , which bounds the expected time-to-goal of the optimal policy from any state. Furthermore, we illustrate various cases (e.g., positive costs, or general costs when an order-accurate estimate of T_{\star} is available) where the regret only contains a logarithmic dependence on T_{\star} , thus yielding the first (nearly) horizon-free regret bound beyond the finite-horizon MDP setting.

[Nested Counterfactual Identification from Arbitrary Surrogate Experiments](#)

- Juan Correa, Sanghack Lee, Elias Bareinboim
- abstract@[open-review](#): The Ladder of Causation describes three qualitatively different types of activities an agent may be interested in engaging in, namely, seeing (observational), doing (interventional), and imagining (counterfactual) (Pearl and Mackenzie, 2018). The inferential challenge imposed by the causal hierarchy is that data is collected by an agent observing or intervening in a system (layers 1 and 2), while its goal may be to understand what would have happened had it taken a different course of action, contrary to what factually ended up happening (layer 3). While there exists a solid understanding of the conditions under which cross-layer inferences are allowed from observations to interventions, the results are somewhat scarcer when targeting counterfactual quantities. In this paper, we study the identification of nested counterfactuals from an arbitrary combination of observations and experiments. Specifically, building on a more explicit definition of nested counterfactuals, we prove the counterfactual unnesting theorem (CUT), which allows one to map arbitrary nested counterfactuals to unnested ones. For instance, applications in mediation and fairness analysis usually evoke notions of direct, indirect, and spurious effects, which naturally require nesting. Second, we introduce a sufficient and necessary graphical condition for counterfactual identification from an arbitrary combination of observational and experimental distributions. Lastly, we develop an efficient and complete algorithm for identifying nested counterfactuals; failure of the algorithm returning an expression for a query implies it is not identifiable.

[Sim and Real: Better Together](#)

- Shirli Di-Castro, Dotan Di Castro, Shie Mannor
- abstract@[open-review](#): Simulation is used extensively in autonomous systems, particularly in robotic manipulation. By far, the most common approach is to train a controller in simulation, and then use it as an initial starting point for the real system. We demonstrate how to learn simultaneously from both simulation and interaction with the real environment. We propose an algorithm for balancing the large number of samples from the high throughput but less accurate simulation and the low-throughput, high-fidelity and costly samples from the real environment. We achieve that by maintaining a replay buffer for each environment the agent interacts with. We analyze such multi-environment interaction theoretically, and provide convergence properties, through a novel theoretical replay buffer analysis. We demonstrate the efficacy of our method on a sim-to-real environment.

Trustworthy Multimodal Regression with Mixture of Normal-inverse Gamma Distributions

- Huan Ma, Zongbo Han, Changqing Zhang, Huazhu Fu, Joey Tianyi Zhou, Qinghua Hu
- abstract@[open-review](#): Multimodal regression is a fundamental task, which integrates the information from different sources to improve the performance of follow-up applications. However, existing methods mainly focus on improving the performance and often ignore the confidence of prediction for diverse situations. In this study, we are devoted to trustworthy multimodal regression which is critical in cost-sensitive domains. To this end, we introduce a novel Mixture of Normal-Inverse Gamma distributions (MoNIG) algorithm, which efficiently estimates uncertainty in principle for adaptive integration of different modalities and produces a trustworthy regression result. Our model can be dynamically aware of uncertainty for each modality, and also robust for corrupted modalities. Furthermore, the proposed MoNIG ensures explicitly representation of (modality-specific/global) epistemic and aleatoric uncertainties, respectively. Experimental results on both synthetic and different real-world data demonstrate the effectiveness and trustworthiness of our method on various multimodal regression tasks (e.g., temperature prediction for superconductivity, relative location prediction for CT slices, and multimodal sentiment analysis).

An Empirical Study of Adder Neural Networks for Object Detection

- Xinghao Chen, Chang Xu, Minjing Dong, Chunjing XU, Yunhe Wang
- abstract@[open-review](#): Adder neural networks (AdderNets) have shown impressive performance on image classification with only addition operations, which are more energy efficient than traditional convolutional neural networks built with multiplications. Compared with classification, there is a strong demand on reducing the energy consumption of modern object detectors via AdderNets for real-world applications such as autonomous driving and face detection. In this paper, we present an empirical study of AdderNets for object detection. We first reveal that the batch normalization statistics in the pre-trained adder backbone should not be frozen, since the relatively large feature variance of AdderNets. Moreover, we insert more shortcut connections in the neck part and design a new feature fusion architecture for avoiding the sparse features of adder layers. We present extensive ablation studies to explore several design choices of adder detectors. Comparisons with state-of-the-arts are conducted on COCO and PASCAL VOC benchmarks. Specifically, the proposed Adder FCOS achieves a 37.8% AP on the COCO val set, demonstrating comparable performance to that of the convolutional counterpart with an about \$1.4\times\$ energy reduction.

Does Knowledge Distillation Really Work?

- Samuel Stanton, Pavel Izmailov, Polina Kirichenko, Alexander A. Alemi, Andrew G. Wilson
- abstract@[open-review](#): Knowledge distillation is a popular technique for training a small student network to emulate a larger teacher model, such as an ensemble of networks. We show that while knowledge distillation can improve student generalization, it does not typically work as it is commonly understood: there often remains a surprisingly large discrepancy between the predictive distributions of the teacher and the student, even in cases when the student has the capacity to perfectly match the teacher. We identify difficulties in optimization as a key reason for why the student is unable to match the teacher. We also show how the details of the dataset used for distillation play a role in how closely the student matches the teacher --- and that more closely matching the teacher paradoxically does not always lead to better student generalization.

Teachable Reinforcement Learning via Advice Distillation

- Olivia Watkins, Abhishek Gupta, Trevor Darrell, Pieter Abbeel, Jacob Andreas
- abstract@[open-review](#): Training automated agents to complete complex tasks in interactive environments is challenging: reinforcement learning requires careful hand-engineering of reward functions, imitation learning requires specialized infrastructure and access to a human expert, and learning from intermediate forms of supervision (like binary preferences) is time-consuming and extracts little information from each human intervention. Can we overcome these challenges by building agents that learn from rich, interactive feedback instead? We propose a new supervision paradigm for interactive learning based on "teachable" decision-making systems that learn from structured advice provided by an external teacher. We begin by formalizing a class of human-in-the-loop decision making problems in which multiple forms of teacher-provided advice are available to a learner. We then describe a simple learning algorithm for these problems that first learns to interpret advice, then learns from advice to complete tasks even in the absence of human supervision. In puzzle-solving, navigation, and locomotion domains, we show that agents that learn from advice can acquire new skills with significantly less human supervision than standard reinforcement learning algorithms and often less than imitation learning.

Antipodes of Label Differential Privacy: PATE and ALIBI

- Mani Malek Esmaeili, Ilya Mironov, Karthik Prasad, Igor Shilov, Florian Tramer
- abstract@[open-review](#): We consider the privacy-preserving machine learning (ML) setting where the trained model must satisfy differential privacy (DP) with respect to the labels of the training examples. We propose two novel approaches based on, respectively, the Laplace mechanism and the PATE framework, and demonstrate their effectiveness on standard benchmarks. While recent work by Ghazi et al. proposed Label DP schemes based on a randomized response mechanism, we argue that additive Laplace noise coupled with Bayesian inference (ALIBI) is a better fit for typical ML tasks. Moreover, we show how to achieve very strong privacy levels in some regimes, with our adaptation of the PATE framework that builds on recent advances in semi-supervised learning. We complement theoretical analysis of our algorithms' privacy guarantees with empirical evaluation of their memorization properties. Our evaluation suggests that comparing different algorithms according to their provable DP guarantees can be misleading and favor a less private algorithm with a tighter analysis. Code for implementation of algorithms and memorization attacks is available from <https://github.com/facebookresearch/labeldpantipodes>.

Visual Search Asymmetry: Deep Nets and Humans Share Similar Inherent Biases

- Shashi Kant Gupta, Mengmi Zhang, CHIA-CHIEN WU, Jeremy Wolfe, Gabriel Kreiman
- abstract@[open-review](#): Visual search is a ubiquitous and often challenging daily task, exemplified by looking for the car keys at home or a friend in a crowd. An intriguing property of some classical search tasks is an asymmetry such that finding a target A among distractors B can be easier than finding B among A. To elucidate the mechanisms responsible for asymmetry in visual search, we propose a computational model that takes a target and a search image as inputs and produces a sequence of eye movements until the target is found. The model integrates eccentricity-dependent visual recognition with target-dependent top-down cues. We compared the model against human behavior in six paradigmatic search tasks that show asymmetry in humans. Without prior exposure to the stimuli or task-specific training, the model provides a plausible mechanism for search asymmetry. We hypothesized that the polarity of search asymmetry arises from experience with the natural environment. We tested this hypothesis by training the model on augmented versions of ImageNet where the biases of natural images were either removed or reversed. The polarity of search asymmetry disappeared or was altered depending on the training protocol. This study highlights how classical perceptual properties can emerge in neural network models, without the need for task-specific training, but rather as a consequence of the statistical properties of the developmental diet fed to the model. All source code and data are publicly available at <https://github.com/kreimanlab/VisualSearchAsymmetry>.

On the Universality of Graph Neural Networks on Large Random Graphs

- Nicolas Keriven, Alberto Bietti, Samuel Vaiter

- abstract@[open-review](#): We study the approximation power of Graph Neural Networks (GNNs) on latent position random graphs. In the large graph limit, GNNs are known to converge to certain “continuous” models known as c-GNNs, which directly enables a study of their approximation power on random graph models. In the absence of input node features however, just as GNNs are limited by the Weisfeiler-Lehman isomorphism test, c-GNNs will be severely limited on simple random graph models. For instance, they will fail to distinguish the communities of a well-separated Stochastic Block Model (SBM) with constant degree function. Thus, we consider recently proposed architectures that augment GNNs with unique node identifiers, referred to as Structural GNNs here (SGNNs). We study the convergence of SGNNs to their continuous counterpart (c-SGNNs) in the large random graph limit, under new conditions on the node identifiers. We then show that c-SGNNs are strictly more powerful than c-GNNs in the continuous limit, and prove their universality on several random graph models of interest, including most SBMs and a large class of random geometric graphs. Our results cover both permutation-invariant and permutation-equivariant architectures.

[Inverse Reinforcement Learning in a Continuous State Space with Formal Guarantees](#)

- Gregory Dexter, Kevin Bello, Jean Honorio
- abstract@[open-review](#): Inverse Reinforcement Learning (IRL) is the problem of finding a reward function which describes observed/known expert behavior. The IRL setting is remarkably useful for automated control, in situations where the reward function is difficult to specify manually or as a means to extract agent preference. In this work, we provide a new IRL algorithm for the continuous state space setting with unknown transition dynamics by modeling the system using a basis of orthonormal functions. Moreover, we provide a proof of correctness and formal guarantees on the sample and time complexity of our algorithm. Finally, we present synthetic experiments to corroborate our theoretical guarantees.

[Adversarial Attacks on Graph Classifiers via Bayesian Optimisation](#)

- Xingchen Wan, Henry Kenlay, Robin Ru, Arno Blaas, Michael A Osborne, Xiaowen Dong
- abstract@[open-review](#): Graph neural networks, a popular class of models effective in a wide range of graph-based learning tasks, have been shown to be vulnerable to adversarial attacks. While the majority of the literature focuses on such vulnerability in node-level classification tasks, little effort has been dedicated to analysing adversarial attacks on graph-level classification, an important problem with numerous real-life applications such as biochemistry and social network analysis. The few existing methods often require unrealistic setups, such as access to internal information of the victim models, or an impractically-large number of queries. We present a novel Bayesian optimisation-based attack method for graph classification models. Our method is black-box, query-efficient and parsimonious with respect to the perturbation applied. We empirically validate the effectiveness and flexibility of the proposed method on a wide range of graph classification tasks involving varying graph properties, constraints and modes of attack. Finally, we analyse common interpretable patterns behind the adversarial samples produced, which may shed further light on the adversarial robustness of graph classification models.

[Regulating algorithmic filtering on social media](#)

- Sarah Cen, Devavrat Shah
- abstract@[open-review](#): By filtering the content that users see, social media platforms have the ability to influence users' perceptions and decisions, from their dining choices to their voting preferences. This influence has drawn scrutiny, with many calling for regulations on filtering algorithms, but designing and enforcing regulations remains challenging. In this work, we examine three questions. First, given a regulation, how would one design an audit to enforce it? Second, does the audit impose a performance cost on the platform? Third, how does the audit affect the content that the platform is incentivized to filter? In response to these questions, we propose a method such that, given a regulation, an auditor can test whether that regulation is met with only black-box access to the filtering algorithm. We then turn to the platform's perspective. The platform's goal is to maximize an objective function while meeting regulation. We find that there are conditions under which the regulation does not place a high performance cost on the platform and, notably, that content diversity can play a key role in aligning the interests of the platform and regulators.

[argmax centroid](#)

- Chengyue Gong, Mao Ye, Qiang Liu
- abstract@[open-review](#): We propose a general method to construct centroid approximation for the distribution of maximum points of a random function (a.k.a. argmax distribution), which finds broad applications in machine learning. Our method optimizes a set of centroid points to compactly approximate the argmax distribution with a simple objective function, without explicitly drawing exact samples from the argmax distribution. Theoretically, the argmax centroid method can be shown to minimize a surrogate of Wasserstein distance between the ground-truth argmax distribution and the centroid approximation under proper conditions. We demonstrate the applicability and effectiveness of our method on a variety of real-world multi-task learning applications, including few-shot image classification, personalized dialogue systems and multi-target domain adaptation.

[Contrastive Learning of Global and Local Video Representations](#)

- shuang ma, Zhaoyang Zeng, Daniel McDuff, Yale Song
- abstract@[open-review](#): Contrastive learning has delivered impressive results for various tasks in the self-supervised regime. However, existing approaches optimize for learning representations specific to downstream scenarios, i.e., global representations suitable for tasks such as classification or local representations for tasks such as detection and localization. While they produce satisfactory results in the intended downstream scenarios, they often fail to generalize to tasks that they were not originally designed for. In this work, we propose to learn video representations that generalize to both the tasks which require global semantic information (e.g., classification) and the tasks that require local fine-grained spatio-temporal information (e.g., localization). We achieve this by optimizing two contrastive objectives that together encourage our model to learn global-local visual information given audio signals. We show that the two objectives mutually improve the generalizability of the learned global-local representations, significantly outperforming their disjointly learned counterparts. We demonstrate our approach on various tasks including action/sound classification, lipreading, deepfake detection, event and sound localization.

[BooVI: Provably Efficient Bootstrapped Value Iteration](#)

- Boyi Liu, Qi Cai, Zhuoran Yang, Zhaoran Wang
- abstract@[open-review](#): Despite the tremendous success of reinforcement learning (RL) with function approximation, efficient exploration remains a significant challenge, both practically and theoretically. In particular, existing theoretically grounded RL algorithms based on upper confidence bounds (UCBs), such as optimistic least-squares value iteration (LSVI), are often incompatible with practically powerful function approximators, such as neural networks. In this paper, we develop a variant of \underline{boo}tstrapped LS\underline{VI}, namely BooVI, which bridges such a gap between practice and theory. Practically, BooVI drives exploration through (re)sampling, making it compatible with general function approximators. Theoretically, BooVI inherits the worst-case $\tilde{O}(\sqrt{d^3 H^3 T})$ -regret of optimistic LSVI in the episodic linear setting. Here d is the feature dimension, H is the episode horizon, and T is the total number of steps.

[Do Wider Neural Networks Really Help Adversarial Robustness?](#)

- Boxi Wu, Jinghui Chen, Deng Cai, Xiaofei He, Quanquan Gu

- abstract@[open-review](#): Adversarial training is a powerful type of defense against adversarial examples. Previous empirical results suggest that adversarial training requires wider networks for better performances. However, it remains elusive how does neural network width affect model robustness. In this paper, we carefully examine the relationship between network width and model robustness. Specifically, we show that the model robustness is closely related to the tradeoff between natural accuracy and perturbation stability, which is controlled by the robust regularization parameter λ . With the same λ , wider networks can achieve better natural accuracy but worse perturbation stability, leading to a potentially worse overall model robustness. To understand the origin of this phenomenon, we further relate the perturbation stability with the network's local Lipschitzness. By leveraging recent results on neural tangent kernels, we theoretically show that wider networks tend to have worse perturbation stability. Our analyses suggest that: 1) the common strategy of first fine-tuning λ on small networks and then directly use it for wide model training could lead to deteriorated model robustness; 2) one needs to properly enlarge λ to unleash the robustness potential of wider models fully. Finally, we propose a new Width Adjusted Regularization (WAR) method that adaptively enlarges λ on wide models and significantly saves the tuning time.

[Exploring the Limits of Out-of-Distribution Detection](#)

- Stanislav Fort, Jie Ren, Balaji Lakshminarayanan
- abstract@[open-review](#): Near out-of-distribution detection (OOD) is a major challenge for deep neural networks. We demonstrate that large-scale pre-trained transformers can significantly improve the state-of-the-art (SOTA) on a range of near OOD tasks across different data modalities. For instance, on CIFAR-100 vs CIFAR-10 OOD detection, we improve the AUROC from 85% (current SOTA) to more than 96% using Vision Transformers pre-trained on ImageNet21k. On a challenging genomics OOD detection benchmark, we improve the AUROC from 66% to 77% using transformer and unsupervised pre-training. To further improve performance, we explore the few-shot outlier exposure setting where a few examples from outlier classes may be available; we show that pre-trained transformers are particularly well-suited for outlier exposure, and that the AUROC of OOD detection on CIFAR-100 vs CIFAR-10 can be improved to 98.7% with just 1 image per OOD class, and 99.46% with 10 images per OOD class. For multi-modal image-text pre-trained transformers such as CLIP, we explore a new way of using just the names of outlier classes as a sole source of information without any accompanying images, and show that this outperforms previous SOTA on standard OOD benchmark tasks.

[ABC: Auxiliary Balanced Classifier for Class-imbalanced Semi-supervised Learning](#)

- Hyuck Lee, Seungjae Shin, Heeyoung Kim
- abstract@[open-review](#): Existing semi-supervised learning (SSL) algorithms typically assume class-balanced datasets, although the class distributions of many real world datasets are imbalanced. In general, classifiers trained on a class-imbalanced dataset are biased toward the majority classes. This issue becomes more problematic for SSL algorithms because they utilize the biased prediction of unlabeled data for training. However, traditional class-imbalanced learning techniques, which are designed for labeled data, cannot be readily combined with SSL algorithms. We propose a scalable class-imbalanced SSL algorithm that can effectively use unlabeled data, while mitigating class imbalance by introducing an auxiliary balanced classifier (ABC) of a single layer, which is attached to a representation layer of an existing SSL algorithm. The ABC is trained with a class-balanced loss of a minibatch, while using high-quality representations learned from all data points in the minibatch using the backbone SSL algorithm to avoid overfitting and information loss. Moreover, we use consistency regularization, a recent SSL technique for utilizing unlabeled data in a modified way, to train the ABC to be balanced among the classes by selecting unlabeled data with the same probability for each class. The proposed algorithm achieves state-of-the-art performance in various class-imbalanced SSL experiments using four benchmark datasets.

[BCD Nets: Scalable Variational Approaches for Bayesian Causal Discovery](#)

- Chris Cundy, Aditya Grover, Stefano Ermon
- abstract@[open-review](#): A structural equation model (SEM) is an effective framework to reason over causal relationships represented via a directed acyclic graph (DAG). Recent advances have enabled effective maximum-likelihood point estimation of DAGs from observational data. However, a point estimate may not accurately capture the uncertainty in inferring the underlying graph in practical scenarios, wherein the true DAG is non-identifiable and/or the observed dataset is limited. We propose Bayesian Causal Discovery Nets (BCD Nets), a variational inference framework for estimating a distribution over DAGs characterizing a linear-Gaussian SEM. Developing a full Bayesian posterior over DAGs is challenging due to the discrete and combinatorial nature of graphs. We analyse key design choices for scalable VI over DAGs, such as 1) the parametrization of DAGs via an expressive variational family, 2) a continuous relaxation that enables low-variance stochastic optimization, and 3) suitable priors over the latent variables. We provide a series of experiments on real and synthetic data showing that BCD Nets outperform maximum-likelihood methods on standard causal discovery metrics such as structural Hamming distance in low data regimes.

[Discovering Dynamic Salient Regions for Spatio-Temporal Graph Neural Networks](#)

- Iulia Duta, Andrei Nicolicioiu, Marius Leordeanu
- abstract@[open-review](#): Graph Neural Networks are perfectly suited to capture latent interactions between various entities in the spatio-temporal domain (e.g. videos). However, when an explicit structure is not available, it is not obvious what atomic elements should be represented as nodes. Current works generally use pre-trained object detectors or fixed, predefined regions to extract graph nodes. Improving upon this, our proposed model learns nodes that dynamically attach to well-delimited salient regions, which are relevant for a higher-level task, without using any object-level supervision. Constructing these localized, adaptive nodes gives our model inductive bias towards object-centric representations and we show that it discovers regions that are well correlated with objects in the video. In extensive ablation studies and experiments on two challenging datasets, we show superior performance to previous graph neural networks models for video classification.

[Information-constrained optimization: can adaptive processing of gradients help?](#)

- Jayadev Acharya, Clement Canonne, Prathamesh Mayekar, Himanshu Tyagi
- abstract@[open-review](#): We revisit first-order optimization under local information constraints such as local privacy, gradient quantization, and computational constraints limiting access to a few coordinates of the gradient. In this setting, the optimization algorithm is not allowed to directly access the complete output of the gradient oracle, but only gets limited information about it subject to the local information constraints. We study the role of adaptivity in processing the gradient output to obtain this limited information from it, and obtain tight or nearly tight bounds for both convex and strongly convex optimization when adaptive gradient processing is allowed.

[Towards Calibrated Model for Long-Tailed Visual Recognition from Prior Perspective](#)

- Zhengzhuo Xu, Zenghao Chai, Chun Yuan
- abstract@[open-review](#): Real-world data universally confronts a severe class-imbalance problem and exhibits a long-tailed distribution, i.e., most labels are associated with limited instances. The naïve models supervised by such datasets would prefer dominant labels, encounter a serious generalization challenge and become poorly calibrated. We propose two novel methods from the prior perspective to alleviate this dilemma. First, we deduce a balance-oriented data augmentation named Uniform Mixup (UniMix) to promote mixup in long-tailed scenarios, which adopts advanced mixing factor and sampler in favor of the minority. Second, motivated by the Bayesian theory, we figure out the Bayes Bias (Bayias), an inherent bias caused by the inconsistency of prior, and compensate it as a modification on standard cross-entropy loss. We further prove that both the proposed methods ensure the classification calibration theoretically and empirically. Extensive experiments verify that our strategies contribute to a better-calibrated model, and their combination achieves state-of-the-art performance on CIFAR-LT, ImageNet-LT, and iNaturalist 2018.

Learning to Draw: Emergent Communication through Sketching

- Daniela Mihai, Jonathon Hare
- abstract@[open-review](#): Evidence that visual communication preceded written language and provided a basis for it goes back to prehistory, in forms such as cave and rock paintings depicting traces of our distant ancestors. Emergent communication research has sought to explore how agents can learn to communicate in order to collaboratively solve tasks. Existing research has focused on language, with a learned communication channel transmitting sequences of discrete tokens between the agents. In this work, we explore a visual communication channel between agents that are allowed to draw with simple strokes. Our agents are parameterised by deep neural networks, and the drawing procedure is differentiable, allowing for end-to-end training. In the framework of a referential communication game, we demonstrate that agents can not only successfully learn to communicate by drawing, but with appropriate inductive biases, can do so in a fashion that humans can interpret. We hope to encourage future research to consider visual communication as a more flexible and directly interpretable alternative of training collaborative agents.

Self-Supervised Learning of Event-Based Optical Flow with Spiking Neural Networks

- Jesse Hagenaars, Federico Paredes-Valles, Guido de Croon
- abstract@[open-review](#): The field of neuromorphic computing promises extremely low-power and low-latency sensing and processing. Challenges in transferring learning algorithms from traditional artificial neural networks (ANNs) to spiking neural networks (SNNs) have so far prevented their application to large-scale, complex regression tasks. Furthermore, realizing a truly asynchronous and fully neuromorphic pipeline that maximally attains the abovementioned benefits involves rethinking the way in which this pipeline takes in and accumulates information. In the case of perception, spikes would be passed as-is and one-by-one between an event camera and an SNN, meaning all temporal integration of information must happen inside the network. In this article, we tackle these two problems. We focus on the complex task of learning to estimate optical flow from event-based camera inputs in a self-supervised manner, and modify the state-of-the-art ANN training pipeline to encode minimal temporal information in its inputs. Moreover, we reformulate the self-supervised loss function for event-based optical flow to improve its convexity. We perform experiments with various types of recurrent ANNs and SNNs using the proposed pipeline. Concerning SNNs, we investigate the effects of elements such as parameter initialization and optimization, surrogate gradient shape, and adaptive neuronal mechanisms. We find that initialization and surrogate gradient width play a crucial part in enabling learning with sparse inputs, while the inclusion of adaptivity and learnable neuronal parameters can improve performance. We show that the performance of the proposed ANNs and SNNs are on par with that of the current state-of-the-art ANNs trained in a self-supervised manner.

On the Value of Infinite Gradients in Variational Autoencoder Models

- Bin Dai, Li Wenliang, David Wipf
- abstract@[open-review](#): A number of recent studies of continuous variational autoencoder (VAE) models have noted, either directly or indirectly, the tendency of various parameter gradients to drift towards infinity during training. Because such gradients could potentially contribute to numerical instabilities, and are often framed as a problematic phenomena to be avoided, it may be tempting to shift to alternative energy functions that guarantee bounded gradients. But it remains an open question: What might the unintended consequences of such a restriction be? To address this issue, we examine how unbounded gradients relate to the regularization of a broad class of autoencoder-based architectures, including VAE models, as applied to data lying on or near a low-dimensional manifold (e.g., natural images). Our main finding is that, if the ultimate goal is to simultaneously avoid over-regularization (high reconstruction errors, sometimes referred to as posterior collapse) and under-regularization (excessive latent dimensions are not pruned from the model), then an autoencoder-based energy function with infinite gradients around optimal representations is provably required per a certain technical sense which we carefully detail. Given that both over- and under-regularization can directly lead to poor generated sample quality or suboptimal feature selection, this result suggests that heuristic modifications to or constraints on the VAE energy function may at times be ill-advised, and large gradients should be accommodated to the extent possible.

Online Robust Reinforcement Learning with Model Uncertainty

- Yue Wang, Shaofeng Zou
- abstract@[open-review](#): Robust reinforcement learning (RL) is to find a policy that optimizes the worst-case performance over an uncertainty set of MDPs. In this paper, we focus on model-free robust RL, where the uncertainty set is defined to be centering at a misspecified MDP that generates samples, and is assumed to be unknown. We develop a sample-based approach to estimate the unknown uncertainty set, and design robust Q-learning algorithm (tabular case) and robust TDC algorithm (function approximation setting), which can be implemented in an online and incremental fashion. For the robust Q-learning algorithm, we prove that it converges to the optimal robust Q function, and for the robust TDC algorithm, we prove that it converges asymptotically to some stationary points. Unlike the results in [Roy et al., 2017], our algorithms do not need any additional conditions on the discount factor to guarantee the convergence. We further characterize the finite-time error bounds of the two algorithms, and show that both the robust Q-learning and robust TDC algorithms converge as fast as their vanilla counterparts (within a constant factor). Our numerical experiments further demonstrate the robustness of our algorithms. Our approach can be readily extended to robustify many other algorithms, e.g., TD, SARSA, and other GTD algorithms.

Neural View Synthesis and Matching for Semi-Supervised Few-Shot Learning of 3D Pose

- Angtian Wang, Shenxiao Mei, Alan L. Yuille, Adam Kortylewski
- abstract@[open-review](#): We study the problem of learning to estimate the 3D object pose from a few labelled examples and a collection of unlabelled data. Our main contribution is a learning framework, neural view synthesis and matching, that can transfer the 3D pose annotation from the labelled to unlabelled images reliably, despite unseen 3D views and nuisance variations such as the object shape, texture, illumination or scene context. In our approach, objects are represented as 3D cuboid meshes composed of feature vectors at each mesh vertex. The model is initialized from a few labelled images and is subsequently used to synthesize feature representations of unseen 3D views. The synthesized views are matched with the feature representations of unlabelled images to generate pseudo-labels of the 3D pose. The pseudo-labelled data is, in turn, used to train the feature extractor such that the features at each mesh vertex are more invariant across varying 3D views of the object. Our model is trained in an EM-type manner alternating between increasing the 3D pose invariance of the feature extractor and annotating unlabelled data through neural view synthesis and matching. We demonstrate the effectiveness of the proposed semi-supervised learning framework for 3D pose estimation on the PASCAL3D+ and KITTI datasets. We find that our approach outperforms all baselines by a wide margin, particularly in an extreme few-shot setting where only 7 annotated images are given. Remarkably, we observe that our model also achieves an exceptional robustness in out-of-distribution scenarios that involve partial occlusion.

Sharp Impossibility Results for Hyper-graph Testing

- Jiashun Jin, Zheng Tracy Ke, Jiajun Liang
- abstract@[open-review](#): In a broad Degree-Corrected Mixed-Membership (DCMM) setting, we test whether a non-uniform hypergraph has only one community or has multiple communities. Since both the null and alternative hypotheses have many unknown parameters, the challenge is, given an alternative, how to identify the null that is hardest to separate from the alternative. We approach this by proposing a degree matching strategy where the main idea is leveraging the theory for tensor scaling to create a least favorable pair of hypotheses. We present a result on standard minimax lower bound theory and a result on Region of Impossibility (which is more informative than the minimax lower bound). We show that our lower bounds are tight by introducing a new test that attains the lower bound up to a logarithmic factor. We also discuss the case where the hypergraphs may have mixed-memberships.

[Evaluating Gradient Inversion Attacks and Defenses in Federated Learning](#)

- Yangsibo Huang, Samyak Gupta, Zhao Song, Kai Li, Sanjeev Arora
- abstract@[open-review](#): Gradient inversion attack (or input recovery from gradient) is an emerging threat to the security and privacy preservation of Federated learning, whereby malicious eavesdroppers or participants in the protocol can recover (partially) the clients' private data. This paper evaluates existing attacks and defenses. We find that some attacks make strong assumptions about the setup. Relaxing such assumptions can substantially weaken these attacks. We then evaluate the benefits of three proposed defense mechanisms against gradient inversion attacks. We show the trade-offs of privacy leakage and data utility of these defense methods, and find that combining them in an appropriate manner makes the attack less effective, even under the original strong assumptions. We also estimate the computation cost of end-to-end recovery of a single image under each evaluated defense. Our findings suggest that the state-of-the-art attacks can currently be defended against with minor data utility loss, as summarized in a list of potential strategies.

[Faster Non-asymptotic Convergence for Double Q-learning](#)

- Lin Zhao, Huaqing Xiong, Yingbin Liang
- abstract@[open-review](#): Double Q-learning (Hasselt, 2010) has gained significant success in practice due to its effectiveness in overcoming the overestimation issue of Q-learning. However, the theoretical understanding of double Q-learning is rather limited. The only existing finite-time analysis was recently established in (Xiong et al. 2020), where the polynomial learning rate adopted in the analysis typically yields a slower convergence rate. This paper tackles the more challenging case of a constant learning rate, and develops new analytical tools that improve the existing convergence rate by orders of magnitude. Specifically, we show that synchronous double Q-learning attains an $\$\\epsilon$ -accurate global optimum with a time complexity of $\$\\tilde{\\Omega}\\left(\\frac{\\ln D}{(1-\\gamma)^7\\epsilon^2}\\right)$, and the asynchronous algorithm achieves a time complexity of $\$\\tilde{\\Omega}\\left(\\frac{L}{(1-\\gamma)^7\\epsilon^2}\\right)$, where D is the cardinality of the state-action space, γ is the discount factor, and L is a parameter related to the sampling strategy for asynchronous double Q-learning. These results improve the existing convergence rate by the order of magnitude in terms of its dependence on all major parameters $(\\epsilon, 1-\\gamma, D, L)$. This paper presents a substantial step toward the full understanding of the fast convergence of double-Q learning.

[Towards Tight Communication Lower Bounds for Distributed Optimisation](#)

- Janne H. Korhonen, Dan Alistarh
- abstract@[open-review](#): We consider a standard distributed optimisation setting where N machines, each holding a d -dimensional function f_i , aim to jointly minimise the sum of the functions $\sum_{i=1}^N f_i(x)$. This problem arises naturally in large-scale distributed optimisation, where a standard solution is to apply variants of (stochastic) gradient descent. We focus on the communication complexity of this problem: our main result provides the first fully unconditional bounds on total number of bits which need to be sent and received by the N machines to solve this problem under point-to-point communication, within a given error-tolerance. Specifically, we show that $\Omega(Nd \log d / N\epsilon)$ total bits need to be communicated between the machines to find an additive ϵ -approximation to the minimum of $\sum_{i=1}^N f_i(x)$. The result holds for both deterministic and randomised algorithms, and, importantly, requires no assumptions on the algorithm structure. The lower bound is tight under certain restrictions on parameter values, and is matched within constant factors for quadratic objectives by a new variant of quantised gradient descent, which we describe and analyse. Our results bring over tools from communication complexity to distributed optimisation, which has potential for further applications.

[Fast Multi-Resolution Transformer Fine-tuning for Extreme Multi-label Text Classification](#)

- Jiong Zhang, Wei-Cheng Chang, Hsiang-Fu Yu, Inderjit Dhillon
- abstract@[open-review](#): Extreme multi-label text classification~(XMC) seeks to find relevant labels from an extreme large label collection for a given text input. Many real-world applications can be formulated as XMC problems, such as recommendation systems, document tagging and semantic search. Recently, transformer based XMC methods, such as X-Transformer and LightXML, have shown significant improvement over other XMC methods. Despite leveraging pre-trained transformer models for text representation, the fine-tuning procedure of transformer models on large label space still has lengthy computational time even with powerful GPUs. In this paper, we propose a novel recursive approach, XR-Transformer to accelerate the procedure through recursively fine-tuning transformer models on a series of multi-resolution objectives related to the original XMC objective function. Empirical results show that XR-Transformer takes significantly less training time compared to other transformer-based XMC models while yielding better state-of-the-art results. In particular, on the public Amazon-3M dataset with 3 million labels, XR-Transformer is not only 20x faster than X-Transformer but also improves the Precision@1 from 51% to 54%.

[HRFormer: High-Resolution Vision Transformer for Dense Predict](#)

- YUHUI YUAN, Rao Fu, Lang Huang, Weihong Lin, Chao Zhang, Xilin Chen, Jingdong Wang
- abstract@[open-review](#): We present a High-Resolution Transformer (HRFormer) that learns high-resolution representations for dense prediction tasks, in contrast to the original Vision Transformer that produces low-resolution representations and has high memory and computational cost. We take advantage of the multi-resolution parallel design introduced in high-resolution convolutional networks (HRNet [45]), along with local-window self-attention that performs self-attention over small non-overlapping image windows [21], for improving the memory and computation efficiency. In addition, we introduce a convolution into the FFN to exchange information across the disconnected image windows. We demonstrate the effectiveness of the HighResolution Transformer on both human pose estimation and semantic segmentation tasks, e.g., HRFormer outperforms Swin transformer [27] by 1.3 AP on COCO pose estimation with 50% fewer parameters and 30% fewer FLOPs. Code is available at: <https://github.com/HRNet/HRFormer>

[Manifold Topology Divergence: a Framework for Comparing Data Manifolds.](#)

- Serguei Barannikov, Ilya Trofimov, Grigorii Sotnikov, Ekaterina Trimbach, Alexander Korotin, Alexander Filippov, Evgeny Burnaev
- abstract@[open-review](#): We propose a framework for comparing data manifolds, aimed, in particular, towards the evaluation of deep generative models. We describe a novel tool, Cross-Barcode(P,Q), that, given a pair of distributions in a high-dimensional space, tracks multiscale topology spacial discrepancies between manifolds on which the distributions are concentrated. Based on the Cross-Barcode, we introduce the Manifold Topology Divergence score (MTop-Divergence) and apply it to assess the performance of deep generative models in various domains: images, 3D-shapes, time-series, and on different datasets: MNIST, Fashion MNIST, SVHN, CIFAR10, FFHQ, market stock data, ShapeNet. We demonstrate that the MTop-Divergence accurately detects various degrees of mode-dropping, intra-mode collapse, mode invention, and image disturbance. Our algorithm scales well (essentially linearly) with the increase of the dimension of the ambient high-dimensional space. It is one of the first TDA-based methodologies that can be applied universally to datasets of different sizes and dimensions, including the ones on which the most recent GANs in the visual domain are trained. The proposed method is domain agnostic and does not rely on pre-trained networks.

[Weak-shot Fine-grained Classification via Similarity Transfer](#)

- Junjie Chen, Li Niu, Liu Liu, Liqing Zhang
- abstract@[open-review](#): Recognizing fine-grained categories remains a challenging task, due to the subtle distinctions among different subordinate categories, which results in the need of abundant annotated samples. To alleviate the data-hungry problem, we consider the problem of learning novel categories from web data with the support of a clean set of base categories, which is referred to as weak-shot learning. In this setting, we propose a method

called SimTrans to transfer pairwise semantic similarity from base categories to novel categories. Specifically, we firstly train a similarity net on clean data, and then leverage the transferred similarity to denoise web training data using two simple yet effective strategies. In addition, we apply adversarial loss on similarity net to enhance the transferability of similarity. Comprehensive experiments demonstrate the effectiveness of our weak-shot setting and our SimTrans method.

[Shape your Space: A Gaussian Mixture Regularization Approach to Deterministic Autoencoders](#)

- Amrutha Saseendran, Kathrin Skubch, Stefan Falkner, Margret Keuper
- abstract@[open-review](#): Variational Autoencoders (VAEs) are powerful probabilistic models to learn representations of complex data distributions. One important limitation of VAEs is the strong prior assumption that latent representations learned by the model follow a simple uni-modal Gaussian distribution. Further, the variational training procedure poses considerable practical challenges. Recently proposed regularized autoencoders offer a deterministic autoencoding framework, that simplifies the original VAE objective and is significantly easier to train. Since these models only provide weak control over the learned latent distribution, they require an ex-post density estimation step to generate samples comparable to those of VAEs. In this paper, we propose a simple and end-to-end trainable deterministic autoencoding framework, that efficiently shapes the latent space of the model during training and utilizes the capacity of expressive multi-modal latent distributions. The proposed training procedure provides direct evidence if the latent distribution adequately captures complex aspects of the encoded data. We show in experiments the expressiveness and sample quality of our model in various challenging continuous and discrete domains. An implementation is available at https://github.com/boschresearch/GMM_DAE.

[An Even More Optimal Stochastic Optimization Algorithm: Minibatching and Interpolation Learning](#)

- Blake E. Woodworth, Nathan Srebro
- abstract@[open-review](#): We present and analyze an algorithm for optimizing smooth and convex or strongly convex objectives using minibatch stochastic gradient estimates. The algorithm is optimal with respect to its dependence on both the minibatch size and minimum expected loss simultaneously. This improves over the optimal method of Lan, which is insensitive to the minimum expected loss; over the optimistic acceleration of Cotter et al., which has suboptimal dependence on the minibatch size; and over the algorithm of Liu and Belkin, which is limited to least squares problems and is also similarly suboptimal. Applied to interpolation learning, the improvement over Cotter et al. and Liu and Belkin translates to a linear, rather than square-root, parallelization speedup.

[Indexed Minimum Empirical Divergence for Unimodal Bandits](#)

- Hassan SABER, Pierre Ménard, Odalric-Ambrym Maillard
- abstract@[open-review](#): We consider a stochastic multi-armed bandit problem specified by a set of one-dimensional family exponential distributions endowed with a unimodal structure. The unimodal structure is of practical relevance for several applications. We introduce IMED-UB, an algorithm that exploits provably optimally the unimodal-structure, by adapting to this setting the Indexed Minimum Empirical Divergence (IMED) algorithm introduced by Honda and Takemura (2015). Owing to our proof technique, we are able to provide a concise finite-time analysis of the IMED-UB algorithm, that is simple and yet yields asymptotic optimality. We finally provide numerical experiments showing that IMED-UB competes favorably with the recently introduced state-of-the-art algorithms.

[SOAT: A Scene- and Object-Aware Transformer for Vision-and-Language Navigation](#)

- Abhinav Moudgil, Arjun Majumdar, Harsh Agrawal, Stefan Lee, Dhruv Batra
- abstract@[open-review](#): Natural language instructions for visual navigation often use scene descriptions (e.g., bedroom) and object references (e.g., green chairs) to provide a breadcrumb trail to a goal location. This work presents a transformer-based vision-and-language navigation (VLN) agent that uses two different visual encoders -- a scene classification network and an object detector -- which produce features that match these two distinct types of visual cues. In our method, scene features contribute high-level contextual information that supports object-level processing. With this design, our model is able to use vision-and-language pretraining (i.e., learning the alignment between images and text from large-scale web data) to substantially improve performance on the Room-to-Room (R2R) and Room-Across-Room (RxR) benchmarks. Specifically, our approach leads to improvements of 1.8% absolute in SPL on R2R and 3.7% absolute in SR on RxR. Our analysis reveals even larger gains for navigation instructions that contain six or more object references, which further suggests that our approach is better able to use object features and align them to references in the instructions.

[A Normative and Biologically Plausible Algorithm for Independent Component Analysis](#)

- Yanis Bahroun, Dmitri Chklovskii, Anirvan Sengupta
- abstract@[open-review](#): The brain effortlessly solves blind source separation (BSS) problems, but the algorithm it uses remains elusive. In signal processing, linear BSS problems are often solved by Independent Component Analysis (ICA). To serve as a model of a biological circuit, the ICA neural network (NN) must satisfy at least the following requirements: 1. The algorithm must operate in the online setting where data samples are streamed one at a time, and the NN computes the sources on the fly without storing any significant fraction of the data in memory. 2. The synaptic weight update is local, i.e., it depends only on the biophysical variables present in the vicinity of a synapse. Here, we propose a novel objective function for ICA from which we derive a biologically plausible NN, including both the neural architecture and the synaptic learning rules. Interestingly, our algorithm relies on modulating synaptic plasticity by the total activity of the output neurons. In the brain, this could be accomplished by neuromodulators, extracellular calcium, local field potential, or nitric oxide.

[Regret Bounds for Gaussian-Process Optimization in Large Domains](#)

- Manuel Wuethrich, Bernhard Schölkopf, Andreas Krause
- abstract@[open-review](#): The goal of this paper is to characterize Gaussian-Process optimization in the setting where the function domain is large relative to the number of admissible function evaluations, i.e., where it is impossible to find the global optimum. We provide upper bounds on the suboptimality (Bayesian simple regret) of the solution found by optimization strategies that are closely related to the widely used expected improvement (EI) and upper confidence bound (UCB) algorithms. These regret bounds illuminate the relationship between the number of evaluations, the domain size (i.e. cardinality of finite domains / Lipschitz constant of the covariance function in continuous domains), and the optimality of the retrieved function value. In particular, we show that even when the number of evaluations is far too small to find the global optimum, we can find nontrivial function values (e.g. values that achieve a certain ratio with the optimal value).

[Deeply Shared Filter Bases for Parameter-Efficient Convolutional Neural Networks](#)

- Woochul Kang, Daeyeon Kim
- abstract@[open-review](#): Modern convolutional neural networks (CNNs) have massive identical convolution blocks, and, hence, recursive sharing of parameters across these blocks has been proposed to reduce the amount of parameters. However, naive sharing of parameters poses many challenges such as limited representational power and the vanishing/exploding gradients problem of recursively shared parameters. In this paper, we present a recursive convolution block design and training method, in which a recursively shareable part, or a filter basis, is separated and learned while effectively avoiding the vanishing/exploding gradients problem during training. We show that the unwieldy vanishing/exploding gradients problem can be controlled by

enforcing the elements of the filter basis orthonormal, and empirically demonstrate that the proposed orthogonality regularization improves the flow of gradients during training. Experimental results on image classification and object detection show that our approach, unlike previous parameter-sharing approaches, does not trade performance to save parameters and consistently outperforms over parameterized counterpart networks. This superior performance demonstrates that the proposed recursive convolution block design and the orthogonality regularization not only prevent performance degradation, but also consistently improve the representation capability while a significant amount of parameters are recursively shared.

[On Optimal Robustness to Adversarial Corruption in Online Decision Problems](#)

- Shinji Ito
- abstract@[open-review](#): This paper considers two fundamental sequential decision-making problems: the problem of prediction with expert advice and the multi-armed bandit problem. We focus on stochastic regimes in which an adversary may corrupt losses, and we investigate what level of robustness can be achieved against adversarial corruption. The main contribution of this paper is to show that optimal robustness can be expressed by a square-root dependency on the amount of corruption. More precisely, we show that two classes of algorithms, anytime Hedge with decreasing learning rate and algorithms with second-order regret bounds, achieve $\$O(\sqrt{\log N}/\Delta + \sqrt{C \log N}/\Delta)$ -regret, where N , Δ , and C represent the number of experts, the gap parameter, and the corruption level, respectively. We further provide a matching lower bound, which means that this regret bound is tight up to a constant factor. For the multi-armed bandit problem, we also provide a nearly-tight lower bound up to a logarithmic factor.

[Directed Spectrum Measures Improve Latent Network Models Of Neural Populations](#)

- Neil Gallagher, Kafui Dzirasa, David Carlson
- abstract@[open-review](#): Systems neuroscience aims to understand how networks of neurons distributed throughout the brain mediate computational tasks. One popular approach to identify those networks is to first calculate measures of neural activity (e.g. power spectra) from multiple brain regions, and then apply a linear factor model to those measures. Critically, despite the established role of directed communication between brain regions in neural computation, measures of directed communication have been rarely utilized in network estimation because they are incompatible with the implicit assumptions of the linear factor model approach. Here, we develop a novel spectral measure of directed communication called the Directed Spectrum (DS). We prove that it is compatible with the implicit assumptions of linear factor models, and we provide a method to estimate the DS. We demonstrate that latent linear factor models of DS measures better capture underlying brain networks in both simulated and real neural recording data compared to available alternatives. Thus, linear factor models of the Directed Spectrum offer neuroscientists a simple and effective way to explicitly model directed communication in networks of neural populations.

[Uncertainty-Based Offline Reinforcement Learning with Diversified Q-Ensemble](#)

- Gaon An, Seungyong Moon, Jang-Hyun Kim, Hyun Oh Song
- abstract@[open-review](#): Offline reinforcement learning (offline RL), which aims to find an optimal policy from a previously collected static dataset, bears algorithmic difficulties due to function approximation errors from out-of-distribution (OOD) data points. To this end, offline RL algorithms adopt either a constraint or a penalty term that explicitly guides the policy to stay close to the given dataset. However, prior methods typically require accurate estimation of the behavior policy or sampling from OOD data points, which themselves can be a non-trivial problem. Moreover, these methods under-utilize the generalization ability of deep neural networks and often fall into suboptimal solutions too close to the given dataset. In this work, we propose an uncertainty-based offline RL method that takes into account the confidence of the Q-value prediction and does not require any estimation or sampling of the data distribution. We show that the clipped Q-learning, a technique widely used in online RL, can be leveraged to successfully penalize OOD data points with high prediction uncertainties. Surprisingly, we find that it is possible to substantially outperform existing offline RL methods on various tasks by simply increasing the number of Q-networks along with the clipped Q-learning. Based on this observation, we propose an ensemble-diversified actor-critic algorithm that reduces the number of required ensemble networks down to a tenth compared to the naive ensemble while achieving state-of-the-art performance on most of the D4RL benchmarks considered.

[Distribution-free inference for regression: discrete, continuous, and in between](#)

- Yonghoon Lee, Rina Barber
- abstract@[open-review](#): In data analysis problems where we are not able to rely on distributional assumptions, what types of inference guarantees can still be obtained? Many popular methods, such as holdout methods, cross-validation methods, and conformal prediction, are able to provide distribution-free guarantees for predictive inference, but the problem of providing inference for the underlying regression function (for example, inference on the conditional mean $\mathbb{E}[Y|X]$) is more challenging. In the setting where the features X are continuously distributed, recent work has established that any confidence interval for $\mathbb{E}[Y|X]$ must have non-vanishing width, even as sample size tends to infinity. At the other extreme, if X takes only a small number of possible values, then inference on $\mathbb{E}[Y|X]$ is trivial to achieve. In this work, we study the problem in settings in between these two extremes. We find that there are several distinct regimes in between the finite setting and the continuous setting, where vanishing-width confidence intervals are achievable if and only if the effective support size of the distribution of X is smaller than the square of the sample size.

[Statistical Inference with M-Estimators on Adaptively Collected Data](#)

- Kelly Zhang, Lucas Janson, Susan Murphy
- abstract@[open-review](#): Bandit algorithms are increasingly used in real-world sequential decision-making problems. Associated with this is an increased desire to be able to use the resulting datasets to answer scientific questions like: Did one type of ad lead to more purchases? In which contexts is a mobile health intervention effective? However, classical statistical approaches fail to provide valid confidence intervals when used with data collected with bandit algorithms. Alternative methods have recently been developed for simple models (e.g., comparison of means). Yet there is a lack of general methods for conducting statistical inference using more complex models on data collected with (contextual) bandit algorithms; for example, current methods cannot be used for valid inference on parameters in a logistic regression model for a binary reward. In this work, we develop theory justifying the use of M-estimators---which includes estimators based on empirical risk minimization as well as maximum likelihood---on data collected with adaptive algorithms, including (contextual) bandit algorithms. Specifically, we show that M-estimators, modified with particular adaptive weights, can be used to construct asymptotically valid confidence regions for a variety of inferential targets.

[NeuroLKH: Combining Deep Learning Model with Lin-Kernighan-Helsgaun Heuristic for Solving the Traveling Salesman Problem](#)

- Liang Xin, Wen Song, Zhiguang Cao, Jie Zhang
- abstract@[open-review](#): We present NeuroLKH, a novel algorithm that combines deep learning with the strong traditional heuristic Lin-Kernighan-Helsgaun (LKH) for solving Traveling Salesman Problem. Specifically, we train a Sparse Graph Network (SGN) with supervised learning for edge scores and unsupervised learning for node penalties, both of which are critical for improving the performance of LKH. Based on the output of SGN, NeuroLKH creates the edge candidate set and transforms edge distances to guide the searching process of LKH. Extensive experiments firmly demonstrate that, by training one model on a wide range of problem sizes, NeuroLKH significantly outperforms LKH and generalizes well to much larger sizes. Also, we show

that NeuroLKH can be applied to other routing problems such as Capacitated Vehicle Routing Problem (CVRP), Pickup and Delivery Problem (PDP), and CVRP with Time Windows (CVRPTW).

[LSH-SMILE: Locality Sensitive Hashing Accelerated Simulation and Learning](#)

- Chonghao Sima, Yexiang Xue
- abstract@[open-review](#): The advancement of deep neural networks over the last decade has enabled progress in scientific knowledge discovery in the form of learning Partial Differential Equations (PDEs) directly from experiment data. Nevertheless, forward simulation and backward learning of large-scale dynamic systems require handling billions of mutually interacting elements, the scale of which overwhelms current computing architectures. We propose Locality Sensitive Hashing Accelerated Simulation and Learning (LSH-SMILE), a unified framework to scale up both forward simulation and backward learning of physics systems. LSH-SMILE takes advantage of (i) the locality of PDE updates, (ii) similar temporal dynamics shared by multiple elements. LSH-SMILE hashes elements with similar dynamics into a single hash bucket and handles their updates at once. This allows LSH-SMILE to scale with respect to the number of non-empty hash buckets, a drastic improvement over conventional approaches. Theoretically, we prove a novel bound on the errors introduced by LSH-SMILE. Experimentally, we demonstrate that LSH-SMILE simulates physics systems at comparable quality with exact approaches, but with way less time and space complexity. Such savings also translate to better learning performance due to LSH-SMILE's ability to propagate gradients over a long duration.

[Meta-learning with an Adaptive Task Scheduler](#)

- Huaxiu Yao, Yu Wang, Ying Wei, Peilin Zhao, Mehrdad Mahdavi, Defu Lian, Chelsea Finn
- abstract@[open-review](#): To benefit the learning of a new task, meta-learning has been proposed to transfer a well-generalized meta-model learned from various meta-training tasks. Existing meta-learning algorithms randomly sample meta-training tasks with a uniform probability, under the assumption that tasks are of equal importance. However, it is likely that tasks are detrimental with noise or imbalanced given a limited number of meta-training tasks. To prevent the meta-model from being corrupted by such detrimental tasks or dominated by tasks in the majority, in this paper, we propose an adaptive task scheduler (ATS) for the meta-training process. In ATS, for the first time, we design a neural scheduler to decide which meta-training tasks to use next by predicting the probability being sampled for each candidate task, and train the scheduler to optimize the generalization capacity of the meta-model to unseen tasks. We identify two meta-model-related factors as the input of the neural scheduler, which characterize the difficulty of a candidate task to the meta-model. Theoretically, we show that a scheduler taking the two factors into account improves the meta-training loss and also the optimization landscape. Under the setting of meta-learning with noise and limited budgets, ATS improves the performance on both miniImageNet and a real-world drug discovery benchmark by up to 13% and 18%, respectively, compared to state-of-the-art task schedulers.

[Neural Active Learning with Performance Guarantees](#)

- Zhilei Wang, Pranjal Awasthi, Christoph Dann, Ayush Sekhari, Claudio Gentile
- abstract@[open-review](#): We investigate the problem of active learning in the streaming setting in non-parametric regimes, where the labels are stochastically generated from a class of functions on which we make no assumptions whatsoever. We rely on recently proposed Neural Tangent Kernel (NTK) approximation tools to construct a suitable neural embedding that determines the feature space the algorithm operates on and the learned model computed atop. Since the shape of the label requesting threshold is tightly related to the complexity of the function to be learned, which is a-priori unknown, we also derive a version of the algorithm which is agnostic to any prior knowledge. This algorithm relies on a regret balancing scheme to solve the resulting online model selection problem, and is computationally efficient. We prove joint guarantees on the cumulative regret and number of requested labels which depend on the complexity of the labeling function at hand. In the linear case, these guarantees recover known minimax results of the generalization error as a function of the label complexity in a standard statistical learning setting.

[A Gradient Method for Multilevel Optimization](#)

- Ryo Sato, Mirai Tanaka, Akiko Takeda
- abstract@[open-review](#): Although application examples of multilevel optimization have already been discussed since the 1990s, the development of solution methods was almost limited to bilevel cases due to the difficulty of the problem. In recent years, in machine learning, Franceschi et al. have proposed a method for solving bilevel optimization problems by replacing their lower-level problems with the \$T\$ steepest descent update equations with some prechosen iteration number \$T\$. In this paper, we have developed a gradient-based algorithm for multilevel optimization with \$n\$ levels based on their idea and proved that our reformulation asymptotically converges to the original multilevel problem. As far as we know, this is one of the first algorithms with some theoretical guarantee for multilevel optimization. Numerical experiments show that a trilevel hyperparameter learning model considering data poisoning produces more stable prediction results than an existing bilevel hyperparameter learning model in noisy data settings.

[Edge Representation Learning with Hypergraphs](#)

- Jaehyeong Jo, Jinheon Baek, Seul Lee, Dongki Kim, Minki Kang, Sung Ju Hwang
- abstract@[open-review](#): Graph neural networks have recently achieved remarkable success in representing graph-structured data, with rapid progress in both the node embedding and graph pooling methods. Yet, they mostly focus on capturing information from the nodes considering their connectivity, and not much work has been done in representing the edges, which are essential components of a graph. However, for tasks such as graph reconstruction and generation, as well as graph classification tasks for which the edges are important for discrimination, accurately representing edges of a given graph is crucial to the success of the graph representation learning. To this end, we propose a novel edge representation learning framework based on Dual Hypergraph Transformation (DHT), which transforms the edges of a graph into the nodes of a hypergraph. This dual hypergraph construction allows us to apply message-passing techniques for node representations to edges. After obtaining edge representations from the hypergraphs, we then cluster or drop edges to obtain holistic graph-level edge representations. We validate our edge representation learning method with hypergraphs on diverse graph datasets for graph representation and generation performance, on which our method largely outperforms existing graph representation learning methods. Moreover, our edge representation learning and pooling method also largely outperforms state-of-the-art graph pooling methods on graph classification, not only because of its accurate edge representation learning, but also due to its lossless compression of the nodes and removal of irrelevant edges for effective message-passing.

[One Question Answering Model for Many Languages with Cross-lingual Dense Passage Retrieval](#)

- Akari Asai, Xinyan Yu, Jungo Kasai, Hanna Hajishirzi
- abstract@[open-review](#): We present Cross-lingual Open-Retrieval Answer Generation (CORA), the first unified many-to-many question answering (QA) model that can answer questions across many languages, even for ones without language-specific annotated data or knowledge sources. We introduce a new dense passage retrieval algorithm that is trained to retrieve documents across languages for a question. Combined with a multilingual autoregressive generation model, CORA answers directly in the target language without any translation or in-language retrieval modules as used in prior work. We propose an iterative training method that automatically extends annotated data available only in high-resource languages to low-resource ones. Our results show that CORA substantially outperforms the previous state of the art on multilingual open QA benchmarks across 26 languages, 9 of which are unseen during training. Our analyses show the significance of cross-lingual retrieval and generation in many languages, particularly under low-resource settings.

LEADS: Learning Dynamical Systems that Generalize Across Environments

- Yuan Yin, Ibrahim Ayed, Emmanuel de Bézenac, Nicolas Baskiotis, Patrick Gallinari
- abstract@[open-review](#): When modeling dynamical systems from real-world data samples, the distribution of data often changes according to the environment in which they are captured, and the dynamics of the system itself vary from one environment to another. Generalizing across environments thus challenges the conventional frameworks. The classical settings suggest either considering data as i.i.d and learning a single model to cover all situations or learning environment-specific models. Both are sub-optimal: the former disregards the discrepancies between environments leading to biased solutions, while the latter does not exploit their potential commonalities and is prone to scarcity problems. We propose LEADS, a novel framework that leverages the commonalities and discrepancies among known environments to improve model generalization. This is achieved with a tailored training formulation aiming at capturing common dynamics within a shared model while additional terms capture environment-specific dynamics. We ground our approach in theory, exhibiting a decrease in sample complexity w.r.t classical alternatives. We show how theory and practice coincides on the simplified case of linear dynamics. Moreover, we instantiate this framework for neural networks and evaluate it experimentally on representative families of nonlinear dynamics. We show that this new setting can exploit knowledge extracted from environment-dependent data and improves generalization for both known and novel environments.

Stochastic: A Framework for General Stochastic Automatic Differentiation

- Emile Krieken, Jakub Tomczak, Annette Ten Teije
- abstract@[open-review](#): Modelers use automatic differentiation (AD) of computation graphs to implement complex Deep Learning models without defining gradient computations. Stochastic AD extends AD to stochastic computation graphs with sampling steps, which arise when modelers handle the intractable expectations common in Reinforcement Learning and Variational Inference. However, current methods for stochastic AD are limited: They are either only applicable to continuous random variables and differentiable functions, or can only use simple but high variance score-function estimators. To overcome these limitations, we introduce Stochastic, a new framework for AD of stochastic computation graphs. Stochastic allows the modeler to choose from a wide variety of gradient estimation methods at each sampling step, to optimally reduce the variance of the gradient estimates. Furthermore, Stochastic is provably unbiased for estimation of any-order gradients, and generalizes variance reduction techniques to higher-order gradient estimates. Finally, we implement Stochastic as a PyTorch library at github.com/HEmile/stochastic.

Concentration inequalities under sub-Gaussian and sub-exponential conditions

- Andreas Maurer, Massimiliano Pontil
- abstract@[open-review](#): We prove analogues of the popular bounded difference inequality (also called McDiarmid's inequality) for functions of independent random variables under sub-gaussian and sub-exponential conditions. Applied to vector-valued concentration and the method of Rademacher complexities these inequalities allow an easy extension of uniform convergence results for PCA and linear regression to the case potentially unbounded input- and output variables.

Variance-Aware Off-Policy Evaluation with Linear Function Approximation

- Yifei Min, Tianhao Wang, Dongruo Zhou, Quanquan Gu
- abstract@[open-review](#): We study the off-policy evaluation (OPE) problem in reinforcement learning with linear function approximation, which aims to estimate the value function of a target policy based on the offline data collected by a behavior policy. We propose to incorporate the variance information of the value function to improve the sample efficiency of OPE. More specifically, for time-inhomogeneous episodic linear Markov decision processes (MDPs), we propose an algorithm, \texttt{VA-OPE}, which uses the estimated variance of the value function to reweight the Bellman residual in Fitted Q-Iteration. We show that our algorithm achieves a tighter error bound than the best-known result. We also provide a fine-grained characterization of the distribution shift between the behavior policy and the target policy. Extensive numerical experiments corroborate our theory.

A Provably Efficient Sample Collection Strategy for Reinforcement Learning

- Jean Tarbouriech, Matteo Pirotta, Michal Valko, Alessandro Lazaric
- abstract@[open-review](#): One of the challenges in online reinforcement learning (RL) is that the agent needs to trade off the exploration of the environment and the exploitation of the samples to optimize its behavior. Whether we optimize for regret, sample complexity, state-space coverage or model estimation, we need to strike a different exploration-exploitation trade-off. In this paper, we propose to tackle the exploration-exploitation problem following a decoupled approach composed of: 1) An "objective-specific" algorithm that (adaptively) prescribes how many samples to collect at which states, as if it has access to a generative model (i.e., a simulator of the environment); 2) An "objective-agnostic" sample collection exploration strategy responsible for generating the prescribed samples as fast as possible. Building on recent methods for exploration in the stochastic shortest path problem, we first provide an algorithm that, given as input the number of samples $b(s,a)$ needed in each state-action pair, requires $\tilde{O}(B D + D^{3/2} S^2 A)$ time steps to collect the $B = \sum_{s,a} b(s,a)$ desired samples, in any unknown communicating MDP with S states, A actions and diameter D . Then we show how this general-purpose exploration algorithm can be paired with "objective-specific" strategies that prescribe the sample requirements to tackle a variety of settings — e.g., model estimation, sparse reward discovery, goal-free cost-free exploration in communicating MDPs — for which we obtain improved or novel sample complexity guarantees.

Improved Regret Bounds for Tracking Experts with Memory

- James Robinson, Mark Herbster
- abstract@[open-review](#): We address the problem of sequential prediction with expert advice in a non-stationary environment with long-term memory guarantees in the sense of Bousquet and Warmuth [4]. We give a linear-time algorithm that improves on the best known regret bound [27]. This algorithm incorporates a relative entropy projection step. This projection is advantageous over previous weight-sharing approaches in that weight updates may come with implicit costs as in for example portfolio optimization. We give an algorithm to compute this projection step in linear time, which may be of independent interest.

Robustness of Graph Neural Networks at Scale

- Simon Geisler, Tobias Schmidt, Hakan Sirin, Daniel Zügner, Aleksandar Bojchevski, Stephan Günnemann
- abstract@[open-review](#): Graph Neural Networks (GNNs) are increasingly important given their popularity and the diversity of applications. Yet, existing studies of their vulnerability to adversarial attacks rely on relatively small graphs. We address this gap and study how to attack and defend GNNs at scale. We propose two sparsity-aware first-order optimization attacks that maintain an efficient representation despite optimizing over a number of parameters which is quadratic in the number of nodes. We show that common surrogate losses are not well-suited for global attacks on GNNs. Our alternatives can double the attack strength. Moreover, to improve GNNs' reliability we design a robust aggregation function, Soft Median, resulting in an effective defense at all scales. We evaluate our attacks and defense with standard GNNs on graphs more than 100 times larger compared to previous work. We even scale one order of magnitude further by extending our techniques to a scalable GNN.

Random Noise Defense Against Query-Based Black-Box Attacks

- Zeyu Qin, Yanbo Fan, Hongyuan Zha, Baoyuan Wu
- abstract@[open-review](#): The query-based black-box attacks have raised serious threats to machine learning models in many real applications. In this work, we study a lightweight defense method, dubbed Random Noise Defense (RND), which adds proper Gaussian noise to each query. We conduct the theoretical analysis about the effectiveness of RND against query-based black-box attacks and the corresponding adaptive attacks. Our theoretical results reveal that the defense performance of RND is determined by the magnitude ratio between the noise induced by RND and the noise added by the attackers for gradient estimation or local search. The large magnitude ratio leads to the stronger defense performance of RND, and it's also critical for mitigating adaptive attacks. Based on our analysis, we further propose to combine RND with a plausible Gaussian augmentation Fine-tuning (RND-GF). It enables RND to add larger noise to each query while maintaining the clean accuracy to obtain a better trade-off between clean accuracy and defense performance. Additionally, RND can be flexibly combined with the existing defense methods to further boost the adversarial robustness, such as adversarial training (AT). Extensive experiments on CIFAR-10 and ImageNet verify our theoretical findings and the effectiveness of RND and RND-GF.

[SADGA: Structure-Aware Dual Graph Aggregation Network for Text-to-SQL](#)

- Ruichu Cai, Jinjie Yuan, Boyan Xu, Zhifeng Hao
- abstract@[open-review](#): The Text-to-SQL task, aiming to translate the natural language of the questions into SQL queries, has drawn much attention recently. One of the most challenging problems of Text-to-SQL is how to generalize the trained model to the unseen database schemas, also known as the cross-domain Text-to-SQL task. The key lies in the generalizability of (i) the encoding method to model the question and the database schema and (ii) the question-schema linking method to learn the mapping between words in the question and tables/columns in the database schema. Focusing on the above two key issues, we propose a \emph{Structure-Aware Dual Graph Aggregation Network} (SADGA) for cross-domain Text-to-SQL. In SADGA, we adopt the graph structure to provide a unified encoding model for both the natural language question and database schema. Based on the proposed unified modeling, we further devise a structure-aware aggregation method to learn the mapping between the question-graph and schema-graph. The structure-aware aggregation method is featured with \emph{Global Graph Linking}, \emph{Local Graph Linking} and \emph{Dual-Graph Aggregation Mechanism}. We not only study the performance of our proposal empirically but also achieved 3rd place on the challenging Text-to-SQL benchmark Spider at the time of writing.

[Near-Optimal Offline Reinforcement Learning via Double Variance Reduction](#)

- Ming Yin, Yu Bai, Yu-Xiang Wang
- abstract@[open-review](#): We consider the problem of offline reinforcement learning (RL) --- a well-motivated setting of RL that aims at policy optimization using only historical data. Despite its wide applicability, theoretical understandings of offline RL, such as its optimal sample complexity, remain largely open even in basic settings such as \emph{tabular} Markov Decision Processes (MDPs). In this paper, we propose \emph{Off-Policy Double Variance Reduction} (OPDVR), a new variance reduction-based algorithm for offline RL. Our main result shows that OPDVR provably identifies an $\$\\epsilon\$$ -optimal policy with $\$\\widetilde{O}(H^2/d_m\\epsilon^2)\$$ episodes of offline data in the finite-horizon \emph{stationary transition} setting, where H is the horizon length and d_m is the minimal marginal state-action distribution induced by the behavior policy. This improves over the best-known upper bound by a factor of H . Moreover, we establish an information-theoretic lower bound of $\$\\Omega(H^2/d_m\\epsilon^2)\$$ which certifies that OPDVR is optimal up to logarithmic factors. Lastly, we show that OPDVR also achieves rate-optimal sample complexity under alternative settings such as the finite-horizon MDPs with non-stationary transitions and the infinite horizon MDPs with discounted rewards.

[Joint Modeling of Visual Objects and Relations for Scene Graph Generation](#)

- Minghao Xu, Meng Qu, Bingbing Ni, Jian Tang
- abstract@[open-review](#): An in-depth scene understanding usually requires recognizing all the objects and their relations in an image, encoded as a scene graph. Most existing approaches for scene graph generation first independently recognize each object and then predict their relations independently. Though these approaches are very efficient, they ignore the dependency between different objects as well as between their relations. In this paper, we propose a principled approach to jointly predict the entire scene graph by fully capturing the dependency between different objects and between their relations. Specifically, we establish a unified conditional random field (CRF) to model the joint distribution of all the objects and their relations in a scene graph. We carefully design the potential functions to enable relational reasoning among different objects according to knowledge graph embedding methods. We further propose an efficient and effective algorithm for inference based on mean-field variational inference, in which we first provide a warm initialization by independently predicting the objects and their relations according to the current model, followed by a few iterations of relational reasoning. Experimental results on both the relationship retrieval and zero-shot relationship retrieval tasks prove the efficiency and efficacy of our proposed approach.

[Going Beyond Linear Transformers with Recurrent Fast Weight Programmers](#)

- Kazuki Irie, Imanol Schlag, Róbert Csordás, Jürgen Schmidhuber
- abstract@[open-review](#): Transformers with linearised attention ("linear Transformers") have demonstrated the practical scalability and effectiveness of outer product-based Fast Weight Programmers (FWPs) from the '90s. However, the original FWP formulation is more general than the one of linear Transformers: a slow neural network (NN) continually reprograms the weights of a fast NN with arbitrary architecture. In existing linear Transformers, both NNs are feedforward and consist of a single layer. Here we explore new variations by adding recurrence to the slow and fast nets. We evaluate our novel recurrent FWPs (RFWPs) on two synthetic algorithmic tasks (code execution and sequential ListOps), Wikitext-103 language models, and on the Atari 2600 2D game environment. Our models exhibit properties of Transformers and RNNs. In the reinforcement learning setting, we report large improvements over LSTM in several Atari games. Our code is public.

[Reinforced Few-Shot Acquisition Function Learning for Bayesian Optimization](#)

- Bing-Jing Hsieh, Ping-Chun Hsieh, Xi Liu
- abstract@[open-review](#): Bayesian optimization (BO) conventionally relies on handcrafted acquisition functions (AFs) to sequentially determine the sample points. However, it has been widely observed in practice that the best-performing AF in terms of regret can vary significantly under different types of black-box functions. It has remained a challenge to design one AF that can attain the best performance over a wide variety of black-box functions. This paper aims to attack this challenge through the perspective of reinforced few-shot AF learning (FSAF). Specifically, we first connect the notion of AFs with Q-functions and view a deep Q-network (DQN) as a surrogate differentiable AF. While it serves as a natural idea to combine DQN and an existing few-shot learning method, we identify that such a direct combination does not perform well due to severe overfitting, which is particularly critical in BO due to the need of a versatile sampling policy. To address this, we present a Bayesian variant of DQN with the following three features: (i) It learns a distribution of Q-networks as AFs based on the Kullback-Leibler regularization framework. This inherently provides the uncertainty required in sampling for BO and mitigates overfitting. (ii) For the prior of the Bayesian DQN, we propose to use a demo policy induced by an off-the-shelf AF for better training stability. (iii) On the meta-level, we leverage the meta-loss of Bayesian model-agnostic meta-learning, which serves as a natural companion to the proposed FSAF. Moreover, with the proper design of the Q-networks, FSAF is general-purpose in that it is agnostic to the dimension and the cardinality of the input domain. Through extensive experiments, we demonstrate that the FSAF achieves comparable or better regrets than the state-of-the-art benchmarks on a wide variety of synthetic and real-world test functions.

[Forster Decomposition and Learning Halfspaces with Noise](#)

- Ilias Diakonikolas, Daniel Kane, Christos Tzamos
- abstract@[open-review](#): A Forster transform is an operation that turns a multivariate distribution into one with good anti-concentration properties. While a Forster transform does not always exist, we show that any distribution can be efficiently decomposed as a disjoint mixture of few distributions for which a Forster transform exists and can be computed efficiently. As the main application of this result, we obtain the first polynomial-time algorithm for distribution-independent PAC learning of halfspaces in the Massart noise model with strongly polynomial sample complexity, i.e., independent of the bit complexity of the examples. Previous algorithms for this learning problem incurred sample complexity scaling polynomially with the bit complexity, even though such a dependence is not information-theoretically necessary.

[Cortico-cerebellar networks as decoupling neural interfaces](#)

- Joseph Pemberton, Ellen Boven, Richard Apps, Rui Ponte Costa
- abstract@[open-review](#): The brain solves the credit assignment problem remarkably well. For credit to be assigned across neural networks they must, in principle, wait for specific neural computations to finish. How the brain deals with this inherent locking problem has remained unclear. Deep learning methods suffer from similar locking constraints both on the forward and feedback phase. Recently, decoupled neural interfaces (DNIs) were introduced as a solution to the forward and feedback locking problems in deep networks. Here we propose that a specialised brain region, the cerebellum, helps the cerebral cortex solve similar locking problems akin to DNIs. To demonstrate the potential of this framework we introduce a systems-level model in which a recurrent cortical network receives online temporal feedback predictions from a cerebellar module. We test this cortico-cerebellar recurrent neural network (ccRNN) model on a number of sensorimotor (line and digit drawing) and cognitive tasks (pattern recognition and caption generation) that have been shown to be cerebellar-dependent. In all tasks, we observe that ccRNNs facilitates learning while reducing ataxia-like behaviours, consistent with classical experimental observations. Moreover, our model also explains recent behavioural and neuronal observations while making several testable predictions across multiple levels. Overall, our work offers a novel perspective on the cerebellum as a brain-wide decoupling machine for efficient credit assignment and opens a new avenue between deep learning and neuroscience.

[To The Point: Correspondence-driven monocular 3D category reconstruction](#)

- Filippos Kokkinos, Iasonas Kokkinos
- abstract@[open-review](#): We present To The Point (TTP), a method for reconstructing 3D objects from a single image using 2D to 3D correspondences given only foreground masks, a category specific template and optionally sparse keypoints for supervision. We recover a 3D shape from a 2D image by first regressing the 2D positions corresponding to the 3D template vertices and then jointly estimating a rigid camera transform and non-rigid template deformation that optimally explain the 2D positions through the 3D shape projection. By relying on correspondences we use a simple per-sample optimization problem to replace CNN-based regression of camera pose and non-rigid deformation and thereby obtain substantially more accurate 3D reconstructions. We treat this optimization as a differentiable layer and train the whole system in an end-to-end manner using geometry-driven losses. We report systematic quantitative improvements on multiple categories and provide qualitative results comprising diverse shape, poses and texture prediction examples.

[Proper Value Equivalence](#)

- Christopher Grimm, Andre Barreto, Greg Farquhar, David Silver, Satinder Singh
- abstract@[open-review](#): One of the main challenges in model-based reinforcement learning (RL) is to decide which aspects of the environment should be modeled. The value-equivalence (VE) principle proposes a simple answer to this question: a model should capture the aspects of the environment that are relevant for value-based planning. Technically, VE distinguishes models based on a set of policies and a set of functions: a model is said to be VE to the environment if the Bellman operators it induces for the policies yield the correct result when applied to the functions. As the number of policies and functions increase, the set of VE models shrinks, eventually collapsing to a single point corresponding to a perfect model. A fundamental question underlying the VE principle is thus how to select the smallest sets of policies and functions that are sufficient for planning. In this paper we take an important step towards answering this question. We start by generalizing the concept of VE to order-\$k\$ counterparts defined with respect to \$k\$ applications of the Bellman operator. This leads to a family of VE classes that increase in size as \$k \rightarrow \infty\$. In the limit, all functions become value functions, and we have a special instantiation of VE which we call proper VE or simply PVE. Unlike VE, the PVE class may contain multiple models even in the limit when all value functions are used. Crucially, all these models are sufficient for planning, meaning that they will yield an optimal policy despite the fact that they may ignore many aspects of the environment. We construct a loss function for learning PVE models and argue that popular algorithms such as MuZero can be understood as minimizing an upper bound for this loss. We leverage this connection to propose a modification to MuZero and show that it can lead to improved performance in practice.

[Challenges and Opportunities in High Dimensional Variational Inference](#)

- Akash Kumar Dhaka, Alejandro Catalina, Manushi Welandawe, Michael R. Andersen, Jonathan Huggins, Aki Vehtari
- abstract@[open-review](#): Current black-box variational inference (BBVI) methods require the user to make numerous design choices – such as the selection of variational objective and approximating family – yet there is little principled guidance on how to do so. We develop a conceptual framework and set of experimental tools to understand the effects of these choices, which we leverage to propose best practices for maximizing posterior approximation accuracy. Our approach is based on studying the pre-asymptotic tail behavior of the density ratios between the joint distribution and the variational approximation, then exploiting insights and tools from the importance sampling literature. Our framework and supporting experiments help to distinguish between the behavior of BBVI methods for approximating low-dimensional versus moderate-to-high-dimensional posteriors. In the latter case, we show that mass-covering variational objectives are difficult to optimize and do not improve accuracy, but flexible variational families can improve accuracy and the effectiveness of importance sampling – at the cost of additional optimization challenges. Therefore, for moderate-to-high-dimensional posteriors we recommend using the (mode-seeking) exclusive KL divergence since it is the easiest to optimize, and improving the variational family or using model parameter transformations to make the posterior and optimal variational approximation more similar. On the other hand, in low-dimensional settings, we show that heavy-tailed variational families and mass-covering divergences are effective and can increase the chances that the approximation can be improved by importance sampling.

[On the Expressivity of Markov Reward](#)

- David Abel, Will Dabney, Anna Harutyunyan, Mark K. Ho, Michael Littman, Doina Precup, Satinder Singh
- abstract@[open-review](#): Reward is the driving force for reinforcement-learning agents. This paper is dedicated to understanding the expressivity of reward as a way to capture tasks that we would want an agent to perform. We frame this study around three new abstract notions of “task” that might be desirable: (1) a set of acceptable behaviors, (2) a partial ordering over behaviors, or (3) a partial ordering over trajectories. Our main results prove that while reward can express many of these tasks, there exist instances of each task type that no Markov reward function can capture. We then provide a set of polynomial-time algorithms that construct a Markov reward function that allows an agent to optimize tasks of each of these three types, and correctly determine when no such reward function exists. We conclude with an empirical study that corroborates and illustrates our theoretical findings.

[One More Step Towards Reality: Cooperative Bandits with Imperfect Communication](#)

- Udari Madhushani, Abhimanyu Dubey, Naomi Leonard, Alex Pentland

- abstract@[open-review](#): The cooperative bandit problem is increasingly becoming relevant due to its applications in large-scale decision-making. However, most research for this problem focuses exclusively on the setting with perfect communication, whereas in most real-world distributed settings, communication is often over stochastic networks, with arbitrary corruptions and delays. In this paper, we study cooperative bandit learning under three typical real-world communication scenarios, namely, (a) message-passing over stochastic time-varying networks, (b) instantaneous reward-sharing over a network with random delays, and (c) message-passing with adversarially corrupted rewards, including byzantine communication. For each of these environments, we propose decentralized algorithms that achieve competitive performance, along with near-optimal guarantees on the incurred group regret as well. Furthermore, in the setting with perfect communication, we present an improved delayed-update algorithm that outperforms the existing state-of-the-art on various network topologies. Finally, we present tight network-dependent minimax lower bounds on the group regret. Our proposed algorithms are straightforward to implement and obtain competitive empirical performance.

[Multi-Agent Reinforcement Learning in Stochastic Networked Systems](#)

- Yiheng Lin, Guannan Qu, Longbo Huang, Adam Wierman
- abstract@[open-review](#): We study multi-agent reinforcement learning (MARL) in a stochastic network of agents. The objective is to find localized policies that maximize the (discounted) global reward. In general, scalability is a challenge in this setting because the size of the global state/action space can be exponential in the number of agents. Scalable algorithms are only known in cases where dependencies are static, fixed and local, e.g., between neighbors in a fixed, time-invariant underlying graph. In this work, we propose a Scalable Actor Critic framework that applies in settings where the dependencies can be non-local and stochastic, and provide a finite-time error bound that shows how the convergence rate depends on the speed of information spread in the network. Additionally, as a byproduct of our analysis, we obtain novel finite-time convergence results for a general stochastic approximation scheme and for temporal difference learning with state aggregation, which apply beyond the setting of MARL in networked systems.

[Neural Scene Flow Prior](#)

- Xueqian Li, Jhony Kaesemel Pontes, Simon Lucey
- abstract@[open-review](#): Before the deep learning revolution, many perception algorithms were based on runtime optimization in conjunction with a strong prior/regularization penalty. A prime example of this in computer vision is optical and scene flow. Supervised learning has largely displaced the need for explicit regularization. Instead, they rely on large amounts of labeled data to capture prior statistics, which are not always readily available for many problems. Although optimization is employed to learn the neural network, at runtime, the weights of this network are frozen. As a result, these learning solutions are domain-specific and do not generalize well to other statistically different scenarios. This paper revisits the scene flow problem that relies predominantly on runtime optimization and strong regularization. A central innovation here is the inclusion of a neural scene flow prior, which utilizes the architecture of neural networks as a new type of implicit regularizer. Unlike learning-based scene flow methods, optimization occurs at runtime, and our approach needs no offline datasets---making it ideal for deployment in new environments such as autonomous driving. We show that an architecture based exclusively on multilayer perceptrons (MLPs) can be used as a scene flow prior. Our method attains competitive---if not better---results on scene flow benchmarks. Also, our neural prior's implicit and continuous scene flow representation allows us to estimate dense long-term correspondences across a sequence of point clouds. The dense motion information is represented by scene flow fields where points can be propagated through time by integrating motion vectors. We demonstrate such a capability by accumulating a sequence of lidar point clouds.

[The future is log-Gaussian: ResNets and their infinite-depth-and-width limit at initialization](#)

- Mufan Li, Mihai Nica, Dan Roy
- abstract@[open-review](#): Theoretical results show that neural networks can be approximated by Gaussian processes in the infinite-width limit. However, for fully connected networks, it has been previously shown that for any fixed network width, $\$n\$$, the Gaussian approximation gets worse as the network depth, $\$d\$$, increases. Given that modern networks are deep, this raises the question of how well modern architectures, like ResNets, are captured by the infinite-width limit. To provide a better approximation, we study ReLU ResNets in the infinite-depth-and-width limit, where $\backslash emph\{both\}$ depth and width tend to infinity as their ratio, $\$d/n\$$, remains constant. In contrast to the Gaussian infinite-width limit, we show theoretically that the network exhibits log-Gaussian behaviour at initialization in the infinite-depth-and-width limit, with parameters depending on the ratio $\$d/n\$$. Using Monte Carlo simulations, we demonstrate that even basic properties of standard ResNet architectures are poorly captured by the Gaussian limit, but remarkably well captured by our log-Gaussian limit. Moreover, our analysis reveals that ReLU ResNets at initialization are hypoactivated: fewer than half of the ReLUs are activated. Additionally, we calculate the interlayer correlations, which have the effect of exponentially increasing the variance of the network output. Based on our analysis, we introduce $\backslash emph\{Balanced ResNets\}$, a simple architecture modification, which eliminates hypoactivation and interlayer correlations and is more amenable to theoretical analysis.

[Grammar-Based Grounded Lexicon Learning](#)

- Jiayuan Mao, Freda Shi, Jiajun Wu, Roger Levy, Josh Tenenbaum
- abstract@[open-review](#): We present Grammar-Based Grounded Language Learning (G2L2), a lexicalist approach toward learning a compositional and grounded meaning representation of language from grounded data, such as paired images and texts. At the core of G2L2 is a collection of lexicon entries, which map each word to a tuple of a syntactic type and a neuro-symbolic semantic program. For example, the word shiny has a syntactic type of adjective; its neuro-symbolic semantic program has the symbolic form $\$\\lambda x.\\textit{filter}(x, \\textbf{SHINY})\$$, where the concept SHINY is associated with a neural network embedding, which will be used to classify shiny objects. Given an input sentence, G2L2 first looks up the lexicon entries associated with each token. It then derives the meaning of the sentence as an executable neuro-symbolic program by composing lexical meanings based on syntax. The recovered meaning programs can be executed on grounded inputs. To facilitate learning in an exponentially-growing compositional space, we introduce a joint parsing and expected execution algorithm, which does local marginalization over derivations to reduce the training time. We evaluate G2L2 on two domains: visual reasoning and language-driven navigation. Results show that G2L2 can generalize from small amounts of data to novel compositions of words.

[Distributed Deep Learning In Open Collaborations](#)

- Michael Diskin, Alexey Bukhtiyarov, Max Ryabinin, Lucile Saulnier, quentin lhoest, Anton Sinitin, Dmitry Popov, Dmitry V. Pyrkin, Maxim Kashirin, Alexander Borzunov, Albert Villanova del Moral, Denis Mazur, Ilia Kobelev, Yacine Jernite, Thomas Wolf, Gennady Pekhimenko
- abstract@[open-review](#): Modern deep learning applications require increasingly more compute to train state-of-the-art models. To address this demand, large corporations and institutions use dedicated High-Performance Computing clusters, whose construction and maintenance are both environmentally costly and well beyond the budget of most organizations. As a result, some research directions become the exclusive domain of a few large industrial and even fewer academic actors. To alleviate this disparity, smaller groups may pool their computational resources and run collaborative experiments that benefit all participants. This paradigm, known as grid- or volunteer computing, has seen successful applications in numerous scientific areas. However, using this approach for machine learning is difficult due to high latency, asymmetric bandwidth, and several challenges unique to volunteer computing. In this work, we carefully analyze these constraints and propose a novel algorithmic framework designed specifically for collaborative training. We demonstrate the effectiveness of our approach for SwAV and ALBERT pretraining in realistic conditions and achieve performance comparable to traditional setups at a fraction of the cost. Finally, we provide a detailed report of successful collaborative language model pretraining with nearly 50 participants.

[Neural Ensemble Search for Uncertainty Estimation and Dataset Shift](#)

- Sheheryar Zaidi, Arber Zela, Thomas Elsken, Chris C Holmes, Frank Hutter, Yee Teh
- abstract@[open-review](#): Ensembles of neural networks achieve superior performance compared to standalone networks in terms of accuracy, uncertainty calibration and robustness to dataset shift. Deep ensembles, a state-of-the-art method for uncertainty estimation, only ensemble random initializations of a fixed architecture. Instead, we propose two methods for automatically constructing ensembles with varying architectures, which implicitly trade-off individual architectures' strengths against the ensemble's diversity and exploit architectural variation as a source of diversity. On a variety of classification tasks and modern architecture search spaces, we show that the resulting ensembles outperform deep ensembles not only in terms of accuracy but also uncertainty calibration and robustness to dataset shift. Our further analysis and ablation studies provide evidence of higher ensemble diversity due to architectural variation, resulting in ensembles that can outperform deep ensembles, even when having weaker average base learners. To foster reproducibility, our code is available: <https://github.com/automl/nes>

[Finding Bipartite Components in Hypergraphs](#)

- Peter Macgregor, He Sun
- abstract@[open-review](#): Hypergraphs are important objects to model ternary or higher-order relations of objects, and have a number of applications in analysing many complex datasets occurring in practice. In this work we study a new heat diffusion process in hypergraphs, and employ this process to design a polynomial-time algorithm that approximately finds bipartite components in a hypergraph. We theoretically prove the performance of our proposed algorithm, and compare it against the previous state-of-the-art through extensive experimental analysis on both synthetic and real-world datasets. We find that our new algorithm consistently and significantly outperforms the previous state-of-the-art across a wide range of hypergraphs.

[Hit and Lead Discovery with Explorative RL and Fragment-based Molecule Generation](#)

- Soojung Yang, Doyeong Hwang, Seul Lee, Seongok Ryu, Sung Ju Hwang
- abstract@[open-review](#): Recently, utilizing reinforcement learning (RL) to generate molecules with desired properties has been highlighted as a promising strategy for drug design. Molecular docking program -- a physical simulation that estimates protein-small molecule binding affinity -- can be an ideal reward scoring function for RL, as it is a straightforward proxy of the therapeutic potential. Still, two imminent challenges exist for this task. First, the models often fail to generate chemically realistic and pharmacochemically acceptable molecules. Second, the docking score optimization is a difficult exploration problem that involves many local optima and less smooth surface with respect to molecular structure. To tackle these challenges, we propose a novel RL framework that generates pharmacochemically acceptable molecules with large docking scores. Our method -- Fragment-based generative RL with Explorative Experience replay for Drug design (FREED) -- constrains the generated molecules to a realistic and qualified chemical space and effectively explores the space to find drugs by coupling our fragment-based generation method and a novel error-prioritized experience replay (PER). We also show that our model performs well on both de novo and scaffold-based schemes. Our model produces molecules of higher quality compared to existing methods while achieving state-of-the-art performance on two of three targets in terms of the docking scores of the generated molecules. We further show with ablation studies that our method, predictive error-PER (FREED(PE)), significantly improves the model performance.

[Proxy Convexity: A Unified Framework for the Analysis of Neural Networks Trained by Gradient Descent](#)

- Spencer Frei, Quanquan Gu
- abstract@[open-review](#): Although the optimization objectives for learning neural networks are highly non-convex, gradient-based methods have been wildly successful at learning neural networks in practice. This juxtaposition has led to a number of recent studies on provable guarantees for neural networks trained by gradient descent. Unfortunately, the techniques in these works are often highly specific to the particular setup in each problem, making it difficult to generalize across different settings. To address this drawback in the literature, we propose a unified non-convex optimization framework for the analysis of neural network training. We introduce the notions of proxy convexity and proxy Polyak-Lojasiewicz (PL) inequalities, which are satisfied if the original objective function induces a proxy objective function that is implicitly minimized when using gradient methods. We show that stochastic gradient descent (SGD) on objectives satisfying proxy convexity or the proxy PL inequality leads to efficient guarantees for proxy objective functions. We further show that many existing guarantees for neural networks trained by gradient descent can be unified through proxy convexity and proxy PL inequalities.

[Covariance-Aware Private Mean Estimation Without Private Covariance Estimation](#)

- Gavin Brown, Marco Gaboardi, Adam Smith, Jonathan Ullman, Lydia Zakynthinou
- abstract@[open-review](#): We present two sample-efficient differentially private mean estimators for d -dimensional (sub)Gaussian distributions with unknown covariance. Informally, given $n \geq d\alpha^2$ samples from such a distribution with mean μ and covariance Σ , our estimators output $\tilde{\mu}$ such that $\|\tilde{\mu} - \mu\| \leq \alpha$, where $\|\cdot\|$ is the Mahalanobis distance. All previous estimators with the same guarantee either require strong a priori bounds on the covariance matrix or require $\Omega(d^{3/2})$ samples. Each of our estimators is based on a simple, general approach to designing differentially private mechanisms, but with novel technical steps to make the estimator private and sample-efficient. Our first estimator samples a point with approximately maximum Tukey depth using the exponential mechanism, but restricted to the set of points of large Tukey depth. Proving that this mechanism is private requires a novel analysis. Our second estimator perturbs the empirical mean of the data set with noise calibrated to the empirical covariance. Only the mean is released, however; the covariance is only used internally. Its sample complexity guarantees hold more generally for subgaussian distributions, albeit with a slightly worse dependence on the privacy parameter. For both estimators, careful preprocessing of the data is required to satisfy differential privacy.

[Label consistency in overfitted generalized \$k\$ -means](#)

- Linfan Zhang, Arash Amini
- abstract@[open-review](#): We provide theoretical guarantees for label consistency in generalized k -means problems, with an emphasis on the overfitted case where the number of clusters used by the algorithm is more than the ground truth. We provide conditions under which the estimated labels are close to a refinement of the true cluster labels. We consider both exact and approximate recovery of the labels. Our results hold for any constant-factor approximation to the k -means problem. The results are also model-free and only based on bounds on the maximum or average distance of the data points to the true cluster centers. These centers themselves are loosely defined and can be taken to be any set of points for which the aforementioned distances can be controlled. We show the usefulness of the results with applications to some manifold clustering problems.

[Open-set Label Noise Can Improve Robustness Against Inherent Label Noise](#)

- Hongxin Wei, Lue Tao, RENCHUNZI XIE, Bo An
- abstract@[open-review](#): Learning with noisy labels is a practically challenging problem in weakly supervised learning. In the existing literature, open-set noises are always considered to be poisonous for generalization, similar to closed-set noises. In this paper, we empirically show that open-set noisy labels can be non-toxic and even benefit the robustness against inherent noisy labels. Inspired by the observations, we propose a simple yet effective regularization by introducing Open-set samples with Dynamic Noisy Labels (ODNL) into training. With ODNL, the extra capacity of the neural network can be largely consumed in a way that does not interfere with learning patterns from clean data. Through the lens of SGD noise, we show that the noises induced by our method are random-direction, conflict-free and biased, which may help the model converge to a flat minimum with superior stability and enforce the model to produce conservative predictions on Out-of-Distribution instances. Extensive experimental results on benchmark datasets with

various types of noisy labels demonstrate that the proposed method not only enhances the performance of many existing robust algorithms but also achieves significant improvement on Out-of-Distribution detection tasks even in the label noise setting.

[The Complexity of Sparse Tensor PCA](#)

- Davin Choo, Tommaso d'Orsi
- abstract@[open-review](#): We study the problem of sparse tensor principal component analysis: given a tensor $\mathbf{Y} = \mathbf{W} + \lambda \mathbf{x}^{\otimes p}$ with \mathbf{W} having i.i.d. Gaussian entries, the goal is to recover the k -sparse unit vector \mathbf{x} . The model captures both sparse PCA (in its Wigner form) and tensor PCA. For the highly sparse regime of $k \leq \sqrt{n}$, we present a family of algorithms that smoothly interpolates between a simple polynomial-time algorithm and the exponential-time exhaustive search algorithm. For any $1 \leq t \leq k$, our algorithms recovers the sparse vector for signal-to-noise ratio $\lambda \geq \tilde{\mathcal{O}}(\sqrt{t} \cdot \log(k/t)^{p/2})$ in time $\tilde{\mathcal{O}}(n^{p+t})$, capturing the state-of-the-art guarantees for the matrix settings (in both the polynomial-time and sub-exponential time regimes). Our results naturally extend to the case of r distinct k -sparse signals with disjoint supports, with guarantees that are independent of the number of spikes. Even in the restricted case of sparse PCA, known algorithms only recover the sparse vectors for $\lambda \geq \tilde{\mathcal{O}}(k \cdot r)$ while our algorithms require $\lambda \geq \tilde{\mathcal{O}}(k)$. Finally, by analyzing the low-degree likelihood ratio, we complement these algorithmic results with rigorous evidence illustrating the trade-offs between signal-to-noise ratio and running time. This lower bound captures the known lower bounds for both sparse PCA and tensor PCA. In this general model, we observe a more intricate three-way trade-off between the number of samples n , the sparsity k , and the tensor power p .

[Learning to Elect](#)

- Cem Anil, Xuchan Bao
- abstract@[open-review](#): Voting systems have a wide range of applications including recommender systems, web search, product design and elections. Limited by the lack of general-purpose analytical tools, it is difficult to hand-engineer desirable voting rules for each use case. For this reason, it is appealing to automatically discover voting rules geared towards each scenario. In this paper, we show that set-input neural network architectures such as Set Transformers, fully-connected graph networks and DeepSets are both theoretically and empirically well-suited for learning voting rules. In particular, we show that these network models can not only mimic a number of existing voting rules to compelling accuracy --- both position-based (such as Plurality and Borda) and comparison-based (such as Kemeny, Copeland and Maximin) --- but also discover near-optimal voting rules that maximize different social welfare functions. Furthermore, the learned voting rules generalize well to different voter utility distributions and election sizes unseen during training.

[KALE Flow: A Relaxed KL Gradient Flow for Probabilities with Disjoint Support](#)

- Pierre Glaser, Michael Arbel, Arthur Gretton
- abstract@[open-review](#): We study the gradient flow for a relaxed approximation to the Kullback-Leibler (KL) divergence between a moving source and a fixed target distribution. This approximation, termed the KALE (KL approximate lower-bound estimator), solves a regularized version of the Fenchel dual problem defining the KL over a restricted class of functions. When using a Reproducing Kernel Hilbert Space (RKHS) to define the function class, we show that the KALE continuously interpolates between the KL and the Maximum Mean Discrepancy (MMD). Like the MMD and other Integral Probability Metrics, the KALE remains well defined for mutually singular distributions. Nonetheless, the KALE inherits from the limiting KL a greater sensitivity to mismatch in the support of the distributions, compared with the MMD. These two properties make the KALE gradient flow particularly well suited when the target distribution is supported on a low-dimensional manifold. Under an assumption of sufficient smoothness of the trajectories, we show the global convergence of the KALE flow. We propose a particle implementation of the flow given initial samples from the source and the target distribution, which we use to empirically confirm the KALE's properties.

[When Is Generalizable Reinforcement Learning Tractable?](#)

- Dhruv Malik, Yuanzhi Li, Pradeep Ravikumar
- abstract@[open-review](#): Agents trained by reinforcement learning (RL) often fail to generalize beyond the environment they were trained in, even when presented with new scenarios that seem similar to the training environment. We study the query complexity required to train RL agents that generalize to multiple environments. Intuitively, tractable generalization is only possible when the environments are similar or close in some sense. To capture this, we introduce Weak Proximity, a natural structural condition that requires the environments to have highly similar transition and reward functions and share a policy providing optimal value. Despite such shared structure, we prove that tractable generalization is impossible in the worst case. This holds even when each individual environment can be efficiently solved to obtain an optimal linear policy, and when the agent possesses a generative model. Our lower bound applies to the more complex task of representation learning for efficient generalization to multiple environments. On the positive side, we introduce Strong Proximity, a strengthened condition which we prove is sufficient for efficient generalization.

[Relational Self-Attention: What's Missing in Attention for Video Understanding](#)

- Manjin Kim, Heeseung Kwon, CHUNYU WANG, Suha Kwak, Minsu Cho
- abstract@[open-review](#): Convolution has been arguably the most important feature transform for modern neural networks, leading to the advance of deep learning. Recent emergence of Transformer networks, which replace convolution layers with self-attention blocks, has revealed the limitation of stationary convolution kernels and opened the door to the era of dynamic feature transforms. The existing dynamic transforms, including self-attention, however, are all limited for video understanding where correspondence relations in space and time, i.e., motion information, are crucial for effective representation. In this work, we introduce a relational feature transform, dubbed the relational self-attention (RSA), that leverages rich structures of spatio-temporal relations in videos by dynamically generating relational kernels and aggregating relational contexts. Our experiments and ablation studies show that the RSA network substantially outperforms convolution and self-attention counterparts, achieving the state of the art on the standard motion-centric benchmarks for video action recognition, such as Something-Something-V1&V2, Diving48, and FineGym.

[Towards Enabling Meta-Learning from Target Models](#)

- Su Lu, Han-Jia Ye, Le Gan, De-Chuan Zhan
- abstract@[open-review](#): Meta-learning can extract an inductive bias from previous learning experience and assist the training of new tasks. It is often realized through optimizing a meta-model with the evaluation loss of task-specific solvers. Most existing algorithms sample non-overlapping \mathcal{S} sets and \mathcal{Q} sets to train and evaluate the solvers respectively due to simplicity (\mathcal{S}/\mathcal{Q} protocol). Different from \mathcal{S}/\mathcal{Q} protocol, we can also evaluate a task-specific solver by comparing it to a target model \mathcal{T} , which is the optimal model for this task or a model that behaves well enough on this task (\mathcal{S}/\mathcal{T} protocol). Although being short of research, \mathcal{S}/\mathcal{T} protocol has unique advantages such as offering more informative supervision, but it is computationally expensive. This paper looks into this special evaluation method and takes a step towards putting it into practice. We find that with a small ratio of tasks armed with target models, classic meta-learning algorithms can be improved a lot without consuming many resources. We empirically verify the effectiveness of \mathcal{S}/\mathcal{T} protocol in a typical application of meta-learning, *i.e.*, few-shot learning. In detail, after constructing target models by fine-tuning the pre-trained network on those hard tasks, we match the task-specific solvers and target models via knowledge distillation.

[A Near-Optimal Algorithm for Debiasing Trained Machine Learning Models](#)

- Ibrahim M. Alabdulmohsin, Mario Lucic
- abstract@[open-review](#): We present a scalable post-processing algorithm for debiasing trained models, including deep neural networks (DNNs), which we prove to be near-optimal by bounding its excess Bayes risk. We empirically validate its advantages on standard benchmark datasets across both classical algorithms as well as modern DNN architectures and demonstrate that it outperforms previous post-processing methods while performing on par with in-processing. In addition, we show that the proposed algorithm is particularly effective for models trained at scale where post-processing is a natural and practical choice.

[GENESIS-V2: Inferring Unordered Object Representations without Iterative Refinement](#)

- Martin Engelcke, Oiwi Parker Jones, Ingmar Posner
- abstract@[open-review](#): Advances in unsupervised learning of object-representations have culminated in the development of a broad range of methods for unsupervised object segmentation and interpretable object-centric scene generation. These methods, however, are limited to simulated and real-world datasets with limited visual complexity. Moreover, object representations are often inferred using RNNs which do not scale well to large images or iterative refinement which avoids imposing an unnatural ordering on objects in an image but requires the a priori initialisation of a fixed number of object representations. In contrast to established paradigms, this work proposes an embedding-based approach in which embeddings of pixels are clustered in a differentiable fashion using a stochastic stick-breaking process. Similar to iterative refinement, this clustering procedure also leads to randomly ordered object representations, but without the need of initialising a fixed number of clusters a priori. This is used to develop a new model, GENESIS-v2, which can infer a variable number of object representations without using RNNs or iterative refinement. We show that GENESIS-v2 performs strongly in comparison to recent baselines in terms of unsupervised image segmentation and object-centric scene generation on established synthetic datasets as well as more complex real-world datasets.

[How Data Augmentation affects Optimization for Linear Regression](#)

- Boris Hanin, Yi Sun
- abstract@[open-review](#): Though data augmentation has rapidly emerged as a key tool for optimization in modern machine learning, a clear picture of how augmentation schedules affect optimization and interact with optimization hyperparameters such as learning rate is nascent. In the spirit of classical convex optimization and recent work on implicit bias, the present work analyzes the effect of augmentation on optimization in the simple convex setting of linear regression with MSE loss. We find joint schedules for learning rate and data augmentation scheme under which augmented gradient descent provably converges and characterize the resulting minimum. Our results apply to arbitrary augmentation schemes, revealing complex interactions between learning rates and augmentations even in the convex setting. Our approach interprets augmented (S)GD as a stochastic optimization method for a time-varying sequence of proxy losses. This gives a unified way to analyze learning rate, batch size, and augmentations ranging from additive noise to random projections. From this perspective, our results, which also give rates of convergence, can be viewed as Monro-Robbins type conditions for augmented (S)GD.

[An Exact Characterization of the Generalization Error for the Gibbs Algorithm](#)

- Gholamali Aminian, Yuheng Bu, Laura Toni, Miguel Rodrigues, Gregory Wornell
- abstract@[open-review](#): Various approaches have been developed to upper bound the generalization error of a supervised learning algorithm. However, existing bounds are often loose and lack of guarantees. As a result, they may fail to characterize the exact generalization ability of a learning algorithm. Our main contribution is an exact characterization of the expected generalization error of the well-known Gibbs algorithm (a.k.a. Gibbs posterior) using symmetrized KL information between the input training samples and the output hypothesis. Our result can be applied to tighten existing expected generalization error and PAC-Bayesian bounds. Our approach is versatile, as it also characterizes the generalization error of the Gibbs algorithm with data-dependent regularizer and that of the Gibbs algorithm in the asymptotic regime, where it converges to the empirical risk minimization algorithm. Of particular relevance, our results highlight the role the symmetrized KL information plays in controlling the generalization error of the Gibbs algorithm.

[Subgaussian and Differentiable Importance Sampling for Off-Policy Evaluation and Learning](#)

- Alberto Maria Metelli, Alessio Russo, Marcello Restelli
- abstract@[open-review](#): Importance Sampling (IS) is a widely used building block for a large variety of off-policy estimation and learning algorithms. However, empirical and theoretical studies have progressively shown that vanilla IS leads to poor estimations whenever the behavioral and target policies are too dissimilar. In this paper, we analyze the theoretical properties of the IS estimator by deriving a novel anticoncentration bound that formalizes the intuition behind its undesired behavior. Then, we propose a new class of IS transformations, based on the notion of power mean. To the best of our knowledge, the resulting estimator is the first to achieve, under certain conditions, two key properties: (i) it displays a subgaussian concentration rate; (ii) it preserves the differentiability in the target distribution. Finally, we provide numerical simulations on both synthetic examples and contextual bandits, in comparison with off-policy evaluation and learning baselines.

[Rethinking gradient sparsification as total error minimization](#)

- Atal Sahu, Aritra Dutta, Ahmed M. Abdelmoniem, Trambak Banerjee, Marco Canini, Panos Kalnis
- abstract@[open-review](#): Gradient compression is a widely-established remedy to tackle the communication bottleneck in distributed training of large deep neural networks (DNNs). Under the error-feedback framework, Top-\$k\$ sparsification, sometimes with \$k\$ as little as 0.1% of the gradient size, enables training to the same model quality as the uncompressed case for a similar iteration count. From the optimization perspective, we find that Top-\$k\$ is the communication-optimal sparsifier given a per-iteration \$k\$ element budget. We argue that to further the benefits of gradient sparsification, especially for DNNs, a different perspective is necessary — one that moves from per-iteration optimality to consider optimality for the entire training. We identify that the total error — the sum of the compression errors for all iterations — encapsulates sparsification throughout training. Then, we propose a communication complexity model that minimizes the total error under a communication budget for the entire training. We find that the hard-threshold sparsifier, a variant of the Top-\$k\$ sparsifier with \$k\$ determined by a constant hard-threshold, is the optimal sparsifier for this model. Motivated by this, we provide convex and non-convex convergence analyses for the hard-threshold sparsifier with error-feedback. We show that hard-threshold has the same asymptotic convergence and linear speedup property as SGD in both the case, and unlike with Top-\$k\$ sparsifier, has no impact due to data-heterogeneity. Our diverse experiments on various DNNs and a logistic regression model demonstrate that the hard-threshold sparsifier is more communication-efficient than Top-\$k\$.

[Approximate optimization of convex functions with outlier noise](#)

- Anindya De, Sanjeev Khanna, Huan Li, MohammadHesam NikpeySalekde
- abstract@[open-review](#): We study the problem of minimizing a convex function given by a zeroth order oracle that is possibly corrupted by \$\{\text{em outlier noise}\}\$. Specifically, we assume the function values at some points of the domain are corrupted arbitrarily by an adversary, with the only restriction being that the total volume of corrupted points is bounded. The goal then is to find a point close to the function's minimizer using access to the corrupted oracle. We first prove a lower bound result showing that, somewhat surprisingly, one cannot hope to approximate the minimizer \$\{\text{em nearly as well}\}\$ as

one might expect, even if one is allowed to ask an unbounded number of queries to the oracle. Complementing this negative result, we then develop an efficient algorithm that outputs a point close to the minimizer of the convex function, where the specific distance matches it exactly, up to constant factors, the distance bound shown in our lower bound result.

[Fair Classification with Adversarial Perturbations](#)

- L. Elisa Celis, Anay Mehrotra, Nisheeth Vishnoi
- abstract@[open-review](#): We study fair classification in the presence of an omniscient adversary that, given an η , is allowed to choose an arbitrary η -fraction of the training samples and arbitrarily perturb their protected attributes. The motivation comes from settings in which protected attributes can be incorrect due to strategic misreporting, malicious actors, or errors in imputation; and prior approaches that make stochastic or independence assumptions on errors may not satisfy their guarantees in this adversarial setting. Our main contribution is an optimization framework to learn fair classifiers in this adversarial setting that comes with provable guarantees on accuracy and fairness. Our framework works with multiple and non-binary protected attributes, is designed for the large class of linear-fractional fairness metrics, and can also handle perturbations besides protected attributes. We prove near-tightness of our framework's guarantees for natural hypothesis classes: no algorithm can have significantly better accuracy and any algorithm with better fairness must have lower accuracy. Empirically, we evaluate the classifiers produced by our framework for statistical rate on real-world and synthetic datasets for a family of adversaries.

[Distributed Saddle-Point Problems Under Data Similarity](#)

- Aleksandr Beznosikov, Gesualdo Scutari, Alexander Rogozin, Alexander Gasnikov
- abstract@[open-review](#): We study solution methods for (strongly-)convex-(strongly)-concave Saddle-Point Problems (SPPs) over networks of two type--master/workers (thus centralized) architectures and mesh (thus decentralized) networks. The local functions at each node are assumed to be similar, due to statistical data similarity or otherwise. We establish lower complexity bounds for a fairly general class of algorithms solving the SPP. We show that a given suboptimality ϵ is achieved over master/workers networks in $\Omega(\Delta \cdot \delta \cdot \mu \cdot \log(1/\epsilon))$ rounds of communications, where δ measures the degree of similarity of the local functions, μ is their strong convexity constant, and Δ is the diameter of the network. The lower communication complexity bound over mesh networks reads $\Omega(\log(1/\rho) \cdot \delta \cdot \mu \cdot \log(1/\epsilon))$, where ρ is the (normalized) eigengap of the gossip matrix used for the communication between neighbouring nodes. We then propose algorithms matching the lower bounds over either types of networks (up to log-factors). We assess the effectiveness of the proposed algorithms on a robust regression problem.

[Combining Latent Space and Structured Kernels for Bayesian Optimization over Combinatorial Spaces](#)

- Aryan Deshwal, Jana Doppa
- abstract@[open-review](#): We consider the problem of optimizing combinatorial spaces (e.g., sequences, trees, and graphs) using expensive black-box function evaluations. For example, optimizing molecules for drug design using physical lab experiments. Bayesian optimization (BO) is an efficient framework for solving such problems by intelligently selecting the inputs with high utility guided by a learned surrogate model. A recent BO approach for combinatorial spaces is through a reduction to BO over continuous spaces by learning a latent representation of structures using deep generative models (DGMs). The selected input from the continuous space is decoded into a discrete structure for performing function evaluation. However, the surrogate model over the latent space only uses the information learned by the DGM, which may not have the desired inductive bias to approximate the target black-box function. To overcome this drawback, this paper proposes a principled approach referred as LADDER. The key idea is to define a novel structure-coupled kernel that explicitly integrates the structural information from decoded structures with the learned latent space representation for better surrogate modeling. Our experiments on real-world benchmarks show that LADDER significantly improves over the BO over latent space method, and performs better or similar to state-of-the-art methods.

[Gradual Domain Adaptation without Indexed Intermediate Domains](#)

- Hong-You Chen, Wei-Lun Chao
- abstract@[open-review](#): The effectiveness of unsupervised domain adaptation degrades when there is a large discrepancy between the source and target domains. Gradual domain adaption (GDA) is one promising way to mitigate such an issue, by leveraging additional unlabeled data that gradually shift from the source to the target. Through sequentially adapting the model along the "indexed" intermediate domains, GDA substantially improves the overall adaptation performance. In practice, however, the extra unlabeled data may not be separated into intermediate domains and indexed properly, limiting the applicability of GDA. In this paper, we investigate how to discover the sequence of intermediate domains when it is not already available. Concretely, we propose a coarse-to-fine framework, which starts with a coarse domain discovery step via progressive domain discriminator training. This coarse domain sequence then undergoes a fine indexing step via a novel cycle-consistency loss, which encourages the next intermediate domain to preserve sufficient discriminative knowledge of the current intermediate domain. The resulting domain sequence can then be used by a GDA algorithm. On benchmark data sets of GDA, we show that our approach, which we name Intermediate DDomain Labeler (IDOL), can lead to comparable or even better adaptation performance compared to the pre-defined domain sequence, making GDA more applicable and robust to the quality of domain sequences. Codes are available at <https://github.com/hongyouchu-IDOL>.

[K-level Reasoning for Zero-Shot Coordination in Hanabi](#)

- Brandon Cui, Hengyuan Hu, Luis Pineda, Jakob Foerster
- abstract@[open-review](#): The standard problem setting in cooperative multi-agent settings is self-play (SP), where the goal is to train a team of agents that works well together. However, optimal SP policies commonly contain arbitrary conventions ("handshakes") and are not compatible with other, independently trained agents or humans. This latter desiderata was recently formalized by [Hu2020-OtherPlay](#) as the zero-shot coordination (ZSC) setting and partially addressed with their Other-Play (OP) algorithm, which showed improved ZSC and human-AI performance in the card game Hanabi. OP assumes access to the symmetries of the environment and prevents agents from breaking these in a mutually incompatible way during training. However, as the authors point out, discovering symmetries for a given environment is a computationally hard problem. Instead, we show that through a simple adaption of k-level reasoning (KLR) [Costa-Gomes2006-K-level](#), synchronously training all levels, we can obtain competitive ZSC and ad-hoc teamplay performance in Hanabi, including when paired with a human-like proxy bot. We also introduce a new method, synchronous-k-level reasoning with a best response (SyKLRBR), which further improves performance on our synchronous KLR by co-training a best response.

[Learning Markov State Abstractions for Deep Reinforcement Learning](#)

- Cameron Allen, Neev Parikh, Omer Gottesman, George Konidaris
- abstract@[open-review](#): A fundamental assumption of reinforcement learning in Markov decision processes (MDPs) is that the relevant decision process is, in fact, Markov. However, when MDPs have rich observations, agents typically learn by way of an abstract state representation, and such representations are not guaranteed to preserve the Markov property. We introduce a novel set of conditions and prove that they are sufficient for learning a Markov abstract state representation. We then describe a practical training procedure that combines inverse model estimation and temporal contrastive learning to learn an abstraction that approximately satisfies these conditions. Our novel training objective is compatible with both online and offline training: it does not require a reward signal, but agents can capitalize on reward information when available. We empirically evaluate our approach on a visual gridworld

domain and a set of continuous control benchmarks. Our approach learns representations that capture the underlying structure of the domain and lead to improved sample efficiency over state-of-the-art deep reinforcement learning with visual features---often matching or exceeding the performance achieved with hand-designed compact state information.

[Towards Deeper Deep Reinforcement Learning with Spectral Normalization](#)

- Nils Bjorck, Carla P. Gomes, Kilian Q. Weinberger
- abstract@[open-review](#): In computer vision and natural language processing, innovations in model architecture that increase model capacity have reliably translated into gains in performance. In stark contrast with this trend, state-of-the-art reinforcement learning (RL) algorithms often use small MLPs, and gains in performance typically originate from algorithmic innovations. It is natural to hypothesize that small datasets in RL necessitate simple models to avoid overfitting; however, this hypothesis is untested. In this paper we investigate how RL agents are affected by exchanging the small MLPs with larger modern networks with skip connections and normalization, focusing specifically on actor-critic algorithms. We empirically verify that naively adopting such architectures leads to instabilities and poor performance, likely contributing to the popularity of simple models in practice. However, we show that dataset size is not the limiting factor, and instead argue that instability from taking gradients through the critic is the culprit. We demonstrate that spectral normalization (SN) can mitigate this issue and enable stable training with large modern architectures. After smoothing with SN, larger models yield significant performance improvements --- suggesting that more ``easy'' gains may be had by focusing on model architectures in addition to algorithmic innovations.

[Functionally Regionalized Knowledge Transfer for Low-resource Drug Discovery](#)

- Huaxiu Yao, Ying Wei, Long-Kai Huang, Ding Xue, Junzhou Huang, Zhenhui (Jessie) Li
- abstract@[open-review](#): More recently, there has been a surge of interest in employing machine learning approaches to expedite the drug discovery process where virtual screening for hit discovery and ADMET prediction for lead optimization play essential roles. One of the main obstacles to the wide success of machine learning approaches in these two tasks is that the number of compounds labeled with activities or ADMET properties is too small to build an effective predictive model. This paper seeks to remedy the problem by transferring the knowledge from previous assays, namely in-vivo experiments, by different laboratories and against various target proteins. To accommodate these wildly different assays and capture the similarity between assays, we propose a functional rationalized meta-learning algorithm FRML for such knowledge transfer. FRML constructs the predictive model with layers of neural sub-networks or so-called functional regions. Building on this, FRML shares an initialization for the weights of the predictive model across all assays, while customizes it to each assay with a region localization network choosing the pertinent regions. The compositionality of the model improves the capacity of generalization to various and even out-of-distribution tasks. Empirical results on both virtual screening and ADMET prediction validate the superiority of FRML over state-of-the-art baselines powered with interpretability in assay relationship.

[Memory-Efficient Approximation Algorithms for Max-k-Cut and Correlation Clustering](#)

- Nimita Shinde, Vishnu Narayanan, James Saunderson
- abstract@[open-review](#): Max-k-Cut and correlation clustering are fundamental graph partitioning problems. For a graph $G=(V,E)$ with n vertices, the methods with the best approximation guarantees for Max-k-Cut and the Max-Agree variant of correlation clustering involve solving SDPs with $\mathcal{O}(n^2)$ constraints and variables. Large-scale instances of SDPs, thus, present a memory bottleneck. In this paper, we develop simple polynomial-time Gaussian sampling-based algorithms for these two problems that use $\mathcal{O}(n+|E|)$ memory and nearly achieve the best existing approximation guarantees. For dense graphs arriving in a stream, we eliminate the dependence on $|E|$ in the storage complexity at the cost of a slightly worse approximation ratio by combining our approach with sparsification.

[Panoptic 3D Scene Reconstruction From a Single RGB Image](#)

- Manuel Dahnert, Ji Hou, Matthias Niessner, Angela Dai
- abstract@[open-review](#): Richly segmented 3D scene reconstructions are an integral basis for many high-level scene understanding tasks, such as for robotics, motion planning, or augmented reality. Existing works in 3D perception from a single RGB image tend to focus on geometric reconstruction only, or geometric reconstruction with semantic segmentation or instance segmentation. Inspired by 2D panoptic segmentation, we propose to unify the tasks of geometric reconstruction, 3D semantic segmentation, and 3D instance segmentation into the task of panoptic 3D scene reconstruction -- from a single RGB image, predicting the complete geometric reconstruction of the scene in the camera frustum of the image, along with semantic and instance segmentations. We propose a new approach for holistic 3D scene understanding from a single RGB image which learns to lift and propagate 2D features from an input image to a 3D volumetric scene representation. Our panoptic 3D reconstruction metric evaluates both geometric reconstruction quality as well as panoptic segmentation. Our experiments demonstrate that our approach for panoptic 3D scene reconstruction outperforms alternative approaches for this task.

[Measuring Generalization with Optimal Transport](#)

- Ching-Yao Chuang, Youssef Mroueh, Kristjan Greenewald, Antonio Torralba, Stefanie Jegelka
- abstract@[open-review](#): Understanding the generalization of deep neural networks is one of the most important tasks in deep learning. Although much progress has been made, theoretical error bounds still often behave disparately from empirical observations. In this work, we develop margin-based generalization bounds, where the margins are normalized with optimal transport costs between independent random subsets sampled from the training distribution. In particular, the optimal transport cost can be interpreted as a generalization of variance which captures the structural properties of the learned feature space. Our bounds robustly predict the generalization error, given training data and network parameters, on large scale datasets. Theoretically, we demonstrate that the concentration and separation of features play crucial roles in generalization, supporting empirical results in the literature.

[Uniform Concentration Bounds toward a Unified Framework for Robust Clustering](#)

- Debolina Paul, Saptarshi Chakraborty, Swagatam Das, Jason Xu
- abstract@[open-review](#): Recent advances in center-based clustering continue to improve upon the drawbacks of Lloyd's celebrated k -means algorithm over \$60\$ years after its introduction. Various methods seek to address poor local minima, sensitivity to outliers, and data that are not well-suited to Euclidean measures of fit, but many are supported largely empirically. Moreover, combining such approaches in a piecemeal manner can result in ad hoc methods, and the limited theoretical results supporting each individual contribution may no longer hold. Toward addressing these issues in a principled way, this paper proposes a cohesive robust framework for center-based clustering under a general class of dissimilarity measures. In particular, we present a rigorous theoretical treatment within a Median-of-Means (MoM) estimation framework, showing that it subsumes several popular k -means variants. In addition to unifying existing methods, we derive uniform concentration bounds that complete their analyses, and bridge these results to the MoM framework via Dudley's chaining arguments. Importantly, we neither require any assumptions on the distribution of the outlying observations nor on the relative number of observations n to features p . We establish strong consistency and an error rate of $O(n^{-1/2})$ under mild conditions, surpassing the best-known results in the literature. The methods are empirically validated thoroughly on real and synthetic datasets.

[Learning Signal-Agnostic Manifolds of Neural Fields](#)

- Yilun Du, Katie Collins, Josh Tenenbaum, Vincent Sitzmann
- abstract@[open-review](#): Deep neural networks have been used widely to learn the latent structure of datasets, across modalities such as images, shapes, and audio signals. However, existing models are generally modality-dependent, requiring custom architectures and objectives to process different classes of signals. We leverage neural fields to capture the underlying structure in image, shape, audio and cross-modal audiovisual domains in a modality-independent manner. We cast our task as one of learning a manifold, where we aim to infer a low-dimensional, locally linear subspace in which our data resides. By enforcing coverage of the manifold, local linearity, and local isometry, our model -- dubbed GEM -- learns to capture the underlying structure of datasets across modalities. We can then travel along linear regions of our manifold to obtain perceptually consistent interpolations between samples, and can further use GEM to recover points on our manifold and glean not only diverse completions of input images, but cross-modal hallucinations of audio or image signals. Finally, we show that by walking across the underlying manifold of GEM, we may generate new samples in our signal domains.

Low-dimensional Structure in the Space of Language Representations is Reflected in Brain Responses

- Richard Antonello, Javier S. Turek, Vy Vo, Alexander Huth
- abstract@[open-review](#): How related are the representations learned by neural language models, translation models, and language tagging tasks? We answer this question by adapting an encoder-decoder transfer learning method from computer vision to investigate the structure among 100 different feature spaces extracted from hidden representations of various networks trained on language tasks. This method reveals a low-dimensional structure where language models and translation models smoothly interpolate between word embeddings, syntactic and semantic tasks, and future word embeddings. We call this low-dimensional structure a language representation embedding because it encodes the relationships between representations needed to process language for a variety of NLP tasks. We find that this representation embedding can predict how well each individual feature space maps to human brain responses to natural language stimuli recorded using fMRI. Additionally, we find that the principal dimension of this structure can be used to create a metric which highlights the brain's natural language processing hierarchy. This suggests that the embedding captures some part of the brain's natural language representation structure.

On the Suboptimality of Thompson Sampling in High Dimensions

- Raymond Zhang, Richard Combes
- abstract@[open-review](#): In this paper we consider Thompson Sampling for combinatorial semi-bandits. We demonstrate that, perhaps surprisingly, Thompson Sampling is sub-optimal for this problem in the sense that its regret scales exponentially in the ambient dimension, and its minimax regret scales almost linearly. This phenomenon occurs under a wide variety of assumptions including both non-linear and linear reward functions in the Bernoulli distribution setting. We also show that including a fixed amount of forced exploration to Thompson Sampling does not alleviate the problem. We complement our theoretical results with numerical results and show that in practice Thompson Sampling indeed can perform very poorly in some high dimension situations.

Learning Debiased and Disentangled Representations for Semantic Segmentation

- Sanghyeok Chu, Dongwan Kim, Bohyung Han
- abstract@[open-review](#): Deep neural networks are susceptible to learn biased models with entangled feature representations, which may lead to subpar performances on various downstream tasks. This is particularly true for under-represented classes, where a lack of diversity in the data exacerbates the tendency. This limitation has been addressed mostly in classification tasks, but there is little study on additional challenges that may appear in more complex dense prediction problems including semantic segmentation. To this end, we propose a model-agnostic and stochastic training scheme for semantic segmentation, which facilitates the learning of debiased and disentangled representations. For each class, we first extract class-specific information from the highly entangled feature map. Then, information related to a randomly sampled class is suppressed by a feature selection process in the feature space. By randomly eliminating certain class information in each training iteration, we effectively reduce feature dependencies among classes, and the model is able to learn more debiased and disentangled feature representations. Models trained with our approach demonstrate strong results on multiple semantic segmentation benchmarks, with especially notable performance gains on under-represented classes.

Diversity Matters When Learning From Ensembles

- Giung Nam, Jongmin Yoon, Yoonho Lee, Juho Lee
- abstract@[open-review](#): Deep ensembles excel in large-scale image classification tasks both in terms of prediction accuracy and calibration. Despite being simple to train, the computation and memory cost of deep ensembles limits their practicability. While some recent works propose to distill an ensemble model into a single model to reduce such costs, there is still a performance gap between the ensemble and distilled models. We propose a simple approach for reducing this gap, i.e., making the distilled performance close to the full ensemble. Our key assumption is that a distilled model should absorb as much function diversity inside the ensemble as possible. We first empirically show that the typical distillation procedure does not effectively transfer such diversity, especially for complex models that achieve near-zero training error. To fix this, we propose a perturbation strategy for distillation that reveals diversity by seeking inputs for which ensemble member outputs disagree. We empirically show that a model distilled with such perturbed samples indeed exhibits enhanced diversity, leading to improved performance.

Locally Valid and Discriminative Prediction Intervals for Deep Learning Models

- Zhen Lin, Shubhendu Trivedi, Jimeng Sun
- abstract@[open-review](#): Crucial for building trust in deep learning models for critical real-world applications is efficient and theoretically sound uncertainty quantification, a task that continues to be challenging. Useful uncertainty information is expected to have two key properties: It should be valid (guaranteeing coverage) and discriminative (more uncertain when the expected risk is high). Moreover, when combined with deep learning (DL) methods, it should be scalable and affect the DL model performance minimally. Most existing Bayesian methods lack frequentist coverage guarantees and usually affect model performance. The few available frequentist methods are rarely discriminative and/or violate coverage guarantees due to unrealistic assumptions. Moreover, many methods are expensive or require substantial modifications to the base neural network. Building upon recent advances in conformal prediction [13, 33] and leveraging the classical idea of kernel regression, we propose Locally Valid and Discriminative prediction intervals (LVD), a simple, efficient, and lightweight method to construct discriminative prediction intervals (PIs) for almost any DL model. With no assumptions on the data distribution, such PIs also offer finite-sample local coverage guarantees (contrasted to the simpler marginal coverage). We empirically verify, using diverse datasets, that besides being the only locally valid method for DL, LVD also exceeds or matches the performance (including coverage rate and prediction accuracy) of existing uncertainty quantification methods, while offering additional benefits in scalability and flexibility.

Personalized Federated Learning With Gaussian Processes

- Idan Achituv, Aviv Shamsian, Aviv Navon, Gal Chechik, Ethan Fetaya
- abstract@[open-review](#): Federated learning aims to learn a global model that performs well on client devices with limited cross-client communication. Personalized federated learning (PFL) further extends this setup to handle data heterogeneity between clients by learning personalized models. A key challenge in this setting is to learn effectively across clients even though each client has unique data that is often limited in size. Here we present pFedGP, a solution to PFL that is based on Gaussian processes (GPs) with deep kernel learning. GPs are highly expressive models that work well in the low data regime due to their Bayesian nature. However, applying GPs to PFL raises multiple challenges. Mainly, GPs performance depends heavily on access to a good kernel function, and learning a kernel requires a large training set. Therefore, we propose learning a shared kernel function across all clients,

parameterized by a neural network, with a personal GP classifier for each client. We further extend pFedGP to include inducing points using two novel methods, the first helps to improve generalization in the low data regime and the second reduces the computational cost. We derive a PAC-Bayes generalization bound on novel clients and empirically show that it gives non-vacuous guarantees. Extensive experiments on standard PFL benchmarks with CIFAR-10, CIFAR-100, and CINIC-10, and on a new setup of learning under input noise show that pFedGP achieves well-calibrated predictions while significantly outperforming baseline methods, reaching up to 21% in accuracy gain.

[Risk Bounds for Over-parameterized Maximum Margin Classification on Sub-Gaussian Mixtures](#)

- Yuan Cao, Quanquan Gu, Mikhail Belkin
- abstract@[open-review](#): Modern machine learning systems such as deep neural networks are often highly over-parameterized so that they can fit the noisy training data exactly, yet they can still achieve small test errors in practice. In this paper, we study this "benign overfitting" phenomenon of the maximum margin classifier for linear classification problems. Specifically, we consider data generated from sub-Gaussian mixtures, and provide a tight risk bound for the maximum margin linear classifier in the over-parameterized setting. Our results precisely characterize the condition under which benign overfitting can occur in linear classification problems, and improve on previous work. They also have direct implications for over-parameterized logistic regression.

[Implicit SVD for Graph Representation Learning](#)

- Sami Abu-El-Haija, Hesham Mostafa, Marcel Nassar, Valentino Crespi, Greg Ver Steeg, Aram Galstyan
- abstract@[open-review](#): Recent improvements in the performance of state-of-the-art (SOTA) methods for Graph Representational Learning (GRL) have come at the cost of significant computational resource requirements for training, e.g., for calculating gradients via backprop over many data epochs. Meanwhile, Singular Value Decomposition (SVD) can find closed-form solutions to convex problems, using merely a handful of epochs. In this paper, we make GRL more computationally tractable for those with modest hardware. We design a framework that computes SVD of *implicitly* defined matrices, and apply this framework to several GRL tasks. For each task, we derive first-order approximation of a SOTA model, where we design (expensive-to-store) matrix \mathbf{M} and train the model, in closed-form, via SVD of \mathbf{M} , without calculating entries of \mathbf{M} . By converging to a unique point in one step, and without calculating gradients, our models show competitive empirical test performance over various graphs such as article citation and biological interaction networks. More importantly, SVD can initialize a deeper model, that is architected to be non-linear almost everywhere, though behaves linearly when its parameters reside on a hyperplane, onto which SVD initializes. The deeper model can then be fine-tuned within only a few epochs. Overall, our algorithm trains hundreds of times faster than state-of-the-art methods, while competing on test empirical performance. We open-source our implementation at: <https://github.com/samihaija/isvd>

[Offline Model-based Adaptable Policy Learning](#)

- Xiong-Hui Chen, Yang Yu, Qingyang Li, Fan-Ming Luo, Zhiwei Qin, Wenjie Shang, Jieping Ye
- abstract@[open-review](#): In reinforcement learning, a promising direction to avoid online trial-and-error costs is learning from an offline dataset. Current offline reinforcement learning methods commonly learn in the policy space constrained to in-support regions by the offline dataset, in order to ensure the robustness of the outcome policies. Such constraints, however, also limit the potential of the outcome policies. In this paper, to release the potential of offline policy learning, we investigate the decision-making problems in out-of-support regions directly and propose offline Model-based Adaptable Policy LEarning (MAPLE). By this approach, instead of learning in in-support regions, we learn an adaptable policy that can adapt its behavior in out-of-support regions when deployed. We conduct experiments on MuJoCo controlling tasks with offline datasets. The results show that the proposed method can make robust decisions in out-of-support regions and achieve better performance than SOTA algorithms.

[Multilingual Pre-training with Universal Dependency Learning](#)

- Kailai Sun, Zuchao Li, Hai Zhao
- abstract@[open-review](#): The pre-trained language model (PrLM) demonstrates domination in downstream natural language processing tasks, in which multilingual PrLM takes advantage of language universality to alleviate the issue of limited resources for low-resource languages. Despite its successes, the performance of multilingual PrLM is still unsatisfactory, when multilingual PrLMs only focus on plain text and ignore obvious universal linguistic structure clues. Existing PrLMs have shown that monolingual linguistic structure knowledge may bring about better performance. Thus we propose a novel multilingual PrLM that supports both explicit universal dependency parsing and implicit language modeling. Syntax in terms of universal dependency parse serves as not only pre-training objective but also learned representation in our model, which brings unprecedented PrLM interpretability and convenience in downstream task use. Our model outperforms two popular multilingual PrLM, multilingual-BERT and XLM-R, on cross-lingual natural language understanding (NLU) benchmarks and linguistic structure parsing datasets, demonstrating the effectiveness and stronger cross-lingual modeling capabilities of our approach.

[Parameter-free HE-friendly Logistic Regression](#)

- Junyoung Byun, Woojin Lee, Jaewook Lee
- abstract@[open-review](#): Privacy in machine learning has been widely recognized as an essential ethical and legal issue, because the data used for machine learning may contain sensitive information. Homomorphic encryption has recently attracted attention as a key solution to preserve privacy in machine learning applications. However, current approaches on the training of encrypted machine learning have relied heavily on hyperparameter selection, which should be avoided owing to the extreme difficulty of conducting validation on encrypted data. In this study, we propose an effective privacy-preserving logistic regression method that is free from the approximation of the sigmoid function and hyperparameter selection. In our framework, a logistic regression model can be transformed into the corresponding ridge regression for the logit function. We provide a theoretical background for our framework by suggesting a new generalization error bound on the encrypted data. Experiments on various real-world data show that our framework achieves better classification results while reducing latency by $\sim 68\%$, compared to the previous models.

[Active clustering for labeling training data](#)

- Quentin Lutz, Elie de Panafieu, Maya Stein, Alex Scott
- abstract@[open-review](#): Gathering training data is a key step of any supervised learning task, and it is both critical and expensive. Critical, because the quantity and quality of the training data has a high impact on the performance of the learned function. Expensive, because most practical cases rely on humans-in-the-loop to label the data. The process of determining the correct labels is much more expensive than comparing two items to see whether they belong to the same class. Thus motivated, we propose a setting for training data gathering where the human experts perform the comparatively cheap task of answering pairwise queries, and the computer groups the items into classes (which can be labeled cheaply at the very end of the process). Given the items, we consider two random models for the classes: one where the set partition they form is drawn uniformly, the other one where each item chooses its class independently following a fixed distribution. In the first model, we characterize the algorithms that minimize the average number of queries required to cluster the items and analyze their complexity. In the second model, we analyze a specific algorithm family, propose as a conjecture that they reach the minimum average number of queries and compare their performance to a random approach. We also propose solutions to handle errors or inconsistencies in the experts' answers.

[Exploring Social Posterior Collapse in Variational Autoencoder for Interaction Modeling](#)

- Chen Tang, Wei Zhan, Masayoshi Tomizuka
- abstract@[open-review](#): Multi-agent behavior modeling and trajectory forecasting are crucial for the safe navigation of autonomous agents in interactive scenarios. Variational Autoencoder (VAE) has been widely applied in multi-agent interaction modeling to generate diverse behavior and learn a low-dimensional representation for interacting systems. However, existing literature did not formally discuss if a VAE-based model can properly encode interaction into its latent space. In this work, we argue that one of the typical formulations of VAEs in multi-agent modeling suffers from an issue we refer to as social posterior collapse, i.e., the model is prone to ignoring historical social context when predicting the future trajectory of an agent. It could cause significant prediction errors and poor generalization performance. We analyze the reason behind this under-explored phenomenon and propose several measures to tackle it. Afterward, we implement the proposed framework and experiment on real-world datasets for multi-agent trajectory prediction. In particular, we propose a novel sparse graph attention message-passing (sparse-GAMP) layer, which helps us detect social posterior collapse in our experiments. In the experiments, we verify that social posterior collapse indeed occurs. Also, the proposed measures are effective in alleviating the issue. As a result, the model attains better generalization performance when historical social context is informative for prediction.

[Ensembling Graph Predictions for AMR Parsing](#)

- Thanh Lam Hoang, Gabriele Picco, Yufang Hou, Young-Suk Lee, Lam Nguyen, Dzung Phan, Vanessa Lopez, Ramon Fernandez Astudillo
- abstract@[open-review](#): In many machine learning tasks, models are trained to predict structure data such as graphs. For example, in natural language processing, it is very common to parse texts into dependency trees or abstract meaning representation (AMR) graphs. On the other hand, ensemble methods combine predictions from multiple models to create a new one that is more robust and accurate than individual predictions. In the literature, there are many ensembling techniques proposed for classification or regression problems, however, ensemble graph prediction has not been studied thoroughly. In this work, we formalize this problem as mining the largest graph that is the most supported by a collection of graph predictions. As the problem is NP-Hard, we propose an efficient heuristic algorithm to approximate the optimal solution. To validate our approach, we carried out experiments in AMR parsing problems. The experimental results demonstrate that the proposed approach can combine the strength of state-of-the-art AMR parsers to create new predictions that are more accurate than any individual models in five standard benchmark datasets.

[On the interplay between data structure and loss function in classification problems](#)

- Stéphane d'Ascoli, Marylou Gabrié, Levent Sagun, Giulio Biroli
- abstract@[open-review](#): One of the central features of modern machine learning models, including deep neural networks, is their generalization ability on structured data in the over-parametrized regime. In this work, we consider an analytically solvable setup to investigate how properties of data impact learning in classification problems, and compare the results obtained for quadratic loss and logistic loss. Using methods from statistical physics, we obtain a precise asymptotic expression for the train and test errors of random feature models trained on a simple model of structured data. The input covariance is built from independent blocks allowing us to tune the saliency of low-dimensional structures and their alignment with respect to the target function. Our results show in particular that in the over-parametrized regime, the impact of data structure on both train and test error curves is greater for logistic loss than for mean-squared loss: the easier the task, the wider the gap in performance between the two losses at the advantage of the logistic. Numerical experiments on MNIST and CIFAR10 confirm our insights.

[Near-optimal Offline and Streaming Algorithms for Learning Non-Linear Dynamical Systems](#)

- Suhas Kowshik, Dheeraj Nagaraj, Prateek Jain, Praneeth Netrapalli
- abstract@[open-review](#): We consider the setting of vector valued non-linear dynamical systems $\mathbf{X}_{\{t+1\}} = \phi(\mathbf{A}^{\{j\}} \mathbf{X}_t) + \eta_t$, where η_t is unbiased noise and $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a known link function that satisfies certain ϵ -expansivity property. The goal is to learn $\mathbf{A}^{\{j\}}$ from a single trajectory $\mathbf{X}_1, \dots, \mathbf{X}_T$ of ϵ -dependent or correlated samples. While the problem is well-studied in the linear case, where ϕ is identity, with optimal error rates even for non-mixing systems, existing results in the non-linear case hold only for mixing systems. In this work, we improve existing results for learning nonlinear systems in a number of ways: a) we provide the first offline algorithm that can learn non-linear dynamical systems without the mixing assumption, b) we significantly improve upon the sample complexity of existing results for mixing systems, c) in the much harder one-pass, streaming setting we study a SGD with Reverse Experience Replay (SGD-RER) method, and demonstrate that for mixing systems, it achieves the same sample complexity as our offline algorithm, d) we justify the expansivity assumption by showing that for the popular ReLU link function --- a non-expansive but easy to learn link function with i.i.d. samples --- any method would require exponentially many samples (with respect to dimension of \mathbf{X}_t) from the dynamical system. We validate our results via simulations and demonstrate that a naive application of SGD can be highly sub-optimal. Indeed, our work demonstrates that for correlated data, specialized methods designed for the dependency structure in data can significantly outperform standard SGD based methods.

[Mixture Proportion Estimation and PU Learning:A Modern Approach](#)

- Saurabh Garg, Yifan Wu, Alexander J. Smola, Sivaraman Balakrishnan, Zachary Lipton
- abstract@[open-review](#): Given only positive examples and unlabeled examples (from both positive and negative classes), we might hope nevertheless to estimate an accurate positive-versus-negative classifier. Formally, this task is broken down into two subtasks: (i) Mixture Proportion Estimation (MPE)---determining the fraction of positive examples in the unlabeled data; and (ii) PU-learning---given such an estimate, learning the desired positive-versus-negative classifier. Unfortunately, classical methods for both problems break down in high-dimensional settings. Meanwhile, recently proposed heuristics lack theoretical coherence and depend precariously on hyperparameter tuning. In this paper, we propose two simple techniques: Best Bin Estimation (BBE) (for MPE); and Conditional Value Ignoring Risk (CVIR), a simple objective for PU-learning. Both methods dominate previous approaches empirically, and for BBE, we establish formal guarantees that hold whenever we can train a model to cleanly separate out a small subset of positive examples. Our final algorithm (TED) n , alternates between the two procedures, significantly improving both our mixture proportion estimator and classifier

[Escape saddle points by a simple gradient-descent based algorithm](#)

- Chenyi Zhang, Tongyang Li
- abstract@[open-review](#): Escaping saddle points is a central research topic in nonconvex optimization. In this paper, we propose a simple gradient-based algorithm such that for a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$, it outputs an ϵ -approximate second-order stationary point in $\tilde{O}(\log n/\epsilon^{1.75})$ iterations. Compared to the previous state-of-the-art algorithms by Jin et al. with $\tilde{O}(\log^4 n/\epsilon^{2})$ or $\tilde{O}(\log^6 n/\epsilon^{1.75})$ iterations, our algorithm is polynomially better in terms of $\log n$ and matches their complexities in terms of $1/\epsilon$. For the stochastic setting, our algorithm outputs an ϵ -approximate second-order stationary point in $\tilde{O}(\log^2 n/\epsilon^4)$ iterations. Technically, our main contribution is an idea of implementing a robust Hessian power method using only gradients, which can find negative curvature near saddle points and achieve the polynomial speedup in $\log n$ compared to the perturbed gradient descent methods. Finally, we also perform numerical experiments that support our results.

[AC/DC: Alternating Compressed/DeCompressed Training of Deep Neural Networks](#)

- Alexandra Peste, Eugenia Iofinova, Adrian Vladu, Dan Alistarh
- abstract@[open-review](#): The increasing computational requirements of deep neural networks (DNNs) have led to significant interest in obtaining DNN models that are sparse, yet accurate. Recent work has investigated the even harder case of sparse training, where the DNN weights are, for as much as

possible, already sparse to reduce computational costs during training. Existing sparse training methods are often empirical and can have lower accuracy relative to the dense baseline. In this paper, we present a general approach called Alternating Compressed/DeCompressed (AC/DC) training of DNNs, demonstrate convergence for a variant of the algorithm, and show that AC/DC outperforms existing sparse training methods in accuracy at similar computational budgets; at high sparsity levels, AC/DC even outperforms existing methods that rely on accurate pre-trained dense models. An important property of AC/DC is that it allows co-training of dense and sparse models, yielding accurate sparse-dense model pairs at the end of the training process. This is useful in practice, where compressed variants may be desirable for deployment in resource-constrained settings without re-doing the entire training flow, and also provides us with insights into the accuracy gap between dense and compressed models.

[HyperSPNs: Compact and Expressive Probabilistic Circuits](#)

- Andy Shih, Dorsa Sadigh, Stefano Ermon
- abstract@[open-review](#): Probabilistic circuits (PCs) are a family of generative models which allows for the computation of exact likelihoods and marginals of its probability distributions. PCs are both expressive and tractable, and serve as popular choices for discrete density estimation tasks. However, large PCs are susceptible to overfitting, and only a few regularization strategies (e.g., dropout, weight-decay) have been explored. We propose HyperSPNs: a new paradigm of generating the mixture weights of large PCs using a small-scale neural network. Our framework can be viewed as a soft weight-sharing strategy, which combines the greater expressiveness of large models with the better generalization and memory-footprint properties of small models. We show the merits of our regularization strategy on two state-of-the-art PC families introduced in recent literature -- RAT-SPNs and EiNETs -- and demonstrate generalization improvements in both models on a suite of density estimation benchmarks in both discrete and continuous domains.

[Scaling Vision with Sparse Mixture of Experts](#)

- Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, Neil Houlsby
- abstract@[open-review](#): Sparsely-gated Mixture of Experts networks (MoEs) have demonstrated excellent scalability in Natural Language Processing. In Computer Vision, however, almost all performant networks are "dense", that is, every input is processed by every parameter. We present a Vision MoE (V-MoE), a sparse version of the Vision Transformer, that is scalable and competitive with the largest dense networks. When applied to image recognition, V-MoE matches the performance of state-of-the-art networks, while requiring as little as half of the compute at inference time. Further, we propose an extension to the routing algorithm that can prioritize subsets of each input across the entire batch, leading to adaptive per-image compute. This allows V-MoE to trade-off performance and compute smoothly at test-time. Finally, we demonstrate the potential of V-MoE to scale vision models, and train a 15B parameter model that attains 90.35% on ImageNet.

[Two-sided fairness in rankings via Lorenz dominance](#)

- Virginie Do, Sam Corbett-Davies, Jamal Atif, Nicolas Usunier
- abstract@[open-review](#): We consider the problem of generating rankings that are fair towards both users and item producers in recommender systems. We address both usual recommendation (e.g., of music or movies) and reciprocal recommendation (e.g., dating). Following concepts of distributive justice in welfare economics, our notion of fairness aims at increasing the utility of the worse-off individuals, which we formalize using the criterion of Lorenz efficiency. It guarantees that rankings are Pareto efficient, and that they maximally redistribute utility from better-off to worse-off, at a given level of overall utility. We propose to generate rankings by maximizing concave welfare functions, and develop an efficient inference procedure based on the Frank-Wolfe algorithm. We prove that unlike existing approaches based on fairness constraints, our approach always produces fair rankings. Our experiments also show that it increases the utility of the worse-off at lower costs in terms of overall utility.

[Stability & Generalisation of Gradient Descent for Shallow Neural Networks without the Neural Tangent Kernel](#)

- Dominic Richards, Ilya Kuzborskij
- abstract@[open-review](#): We revisit on-average algorithmic stability of Gradient Descent (GD) for training overparameterised shallow neural networks and prove new generalisation and excess risk bounds without the Neural Tangent Kernel (NTK) or Polyak-Łojasiewicz (PL) assumptions. In particular, we show oracle type bounds which reveal that the generalisation and excess risk of GD is controlled by an interpolating network with the shortest GD path from initialisation (in a sense, an interpolating network with the smallest relative norm). While this was known for kernelised interpolants, our proof applies directly to networks trained by GD without intermediate kernelisation. At the same time, by relaxing oracle inequalities developed here we recover existing NTK-based risk bounds in a straightforward way, which demonstrates that our analysis is tighter. Finally, unlike most of the NTK-based analyses we focus on regression with label noise and show that GD with early stopping is consistent

[Adversarial Intrinsic Motivation for Reinforcement Learning](#)

- Ishan Durugkar, Mauricio Tec, Scott Niekum, Peter Stone
- abstract@[open-review](#): Learning with an objective to minimize the mismatch with a reference distribution has been shown to be useful for generative modeling and imitation learning. In this paper, we investigate whether one such objective, the Wasserstein-1 distance between a policy's state visitation distribution and a target distribution, can be utilized effectively for reinforcement learning (RL) tasks. Specifically, this paper focuses on goal-conditioned reinforcement learning where the idealized (unachievable) target distribution has full measure at the goal. This paper introduces a quasimetric specific to Markov Decision Processes (MDPs) and uses this quasimetric to estimate the above Wasserstein-1 distance. It further shows that the policy that minimizes this Wasserstein-1 distance is the policy that reaches the goal in as few steps as possible. Our approach, termed Adversarial Intrinsic Motivation (AIM), estimates this Wasserstein-1 distance through its dual objective and uses it to compute a supplemental reward function. Our experiments show that this reward function changes smoothly with respect to transitions in the MDP and directs the agent's exploration to find the goal efficiently. Additionally, we combine AIM with Hindsight Experience Replay (HER) and show that the resulting algorithm accelerates learning significantly on several simulated robotics tasks when compared to other rewards that encourage exploration or accelerate learning.

[Machine Learning for Variance Reduction in Online Experiments](#)

- Yongyi Guo, Dominic Coey, Mikael Konutgan, Wenting Li, Chris Schoener, Matt Goldman
- abstract@[open-review](#): We consider the problem of variance reduction in randomized controlled trials, through the use of covariates correlated with the outcome but independent of the treatment. We propose a machine learning regression-adjusted treatment effect estimator, which we call MLRATE. MLRATE uses machine learning predictors of the outcome to reduce estimator variance. It employs cross-fitting to avoid overfitting biases, and we prove consistency and asymptotic normality under general conditions. MLRATE is robust to poor predictions from the machine learning step: if the predictions are uncorrelated with the outcomes, the estimator performs asymptotically no worse than the standard difference-in-means estimator, while if predictions are highly correlated with outcomes, the efficiency gains are large. In A/A tests, for a set of 48 outcome metrics commonly monitored in Facebook experiments, the estimator has over \$70\%\$ lower variance than the simple difference-in-means estimator, and about \$19\%\$ lower variance than the common univariate procedure which adjusts only for pre-experiment values of the outcome.

[L2ight: Enabling On-Chip Learning for Optical Neural Networks via Efficient in-situ Subspace Optimization](#)

- Jiaqi Gu, Hanqing Zhu, Chenghao Feng, Zixuan Jiang, Ray Chen, David Pan

- abstract@[open-review](#): Silicon-photonics-based optical neural network (ONN) is a promising hardware platform that could represent a paradigm shift in efficient AI with its CMOS-compatibility, flexibility, ultra-low execution latency, and high energy efficiency. In-situ training on the online programmable photonic chips is appealing but still encounters challenging issues in on-chip implementability, scalability, and efficiency. In this work, we propose a closed-loop ONN on-chip learning framework L2ight to enable scalable ONN mapping and efficient in-situ learning. L2ight adopts a three-stage learning flow that first calibrates the complicated photonic circuit states under challenging physical constraints, then performs photonic core mapping via combined analytical solving and zeroth-order optimization. A subspace learning procedure with multi-level sparsity is integrated into L2ight to enable in-situ gradient evaluation and fast adaptation, unleashing the power of optics for real on-chip intelligence. Extensive experiments demonstrate our proposed L2ight outperforms prior ONN training protocols with 3-order-of-magnitude higher scalability and over 30x better efficiency, when benchmarked on various models and learning tasks. This synergistic framework is the first scalable on-chip learning solution that pushes this emerging field from intractable to scalable and further to efficient for next-generation self-learnable photonic neural chips. From a co-design perspective, L2ight also provides essential insights for hardware-restricted unitary subspace optimization and efficient sparse training. We open-source our framework at the link.

[Towards Gradient-based Bilevel Optimization with Non-convex Followers and Beyond](#)

- Risheng Liu, Yaohua Liu, Shangzhi Zeng, Jin Zhang
- abstract@[open-review](#): In recent years, Bi-Level Optimization (BLO) techniques have received extensive attentions from both learning and vision communities. A variety of BLO models in complex and practical tasks are of non-convex follower structure in nature (a.k.a., without Lower-Level Convexity, LLC for short). However, this challenging class of BLOs is lack of developments on both efficient solution strategies and solid theoretical guarantees. In this work, we propose a new algorithmic framework, named Initialization Auxiliary and Pessimistic Trajectory Truncated Gradient Method (IAPTT-GM), to partially address the above issues. In particular, by introducing an auxiliary as initialization to guide the optimization dynamics and designing a pessimistic trajectory truncation operation, we construct a reliable approximate version of the original BLO in the absence of LLC hypothesis. Our theoretical investigations establish the convergence of solutions returned by IAPTT-GM towards those of the original BLO without LLC. As an additional bonus, we also theoretically justify the quality of our IAPTT-GM embedded with Nesterov's accelerated dynamics under LLC. The experimental results confirm both the convergence of our algorithm without LLC, and the theoretical findings under LLC.

[Multi-Facet Clustering Variational Autoencoders](#)

- Fabian Falck, Haoting Zhang, Matthew Willets, George Nicholson, Christopher Yau, Chris C Holmes
- abstract@[open-review](#): Work in deep clustering focuses on finding a single partition of data. However, high-dimensional data, such as images, typically feature multiple interesting characteristics one could cluster over. For example, images of objects against a background could be clustered over the shape of the object and separately by the colour of the background. In this paper, we introduce Multi-Facet Clustering Variational Autoencoders (MFCVAE), a novel class of variational autoencoders with a hierarchy of latent variables, each with a Mixture-of-Gaussians prior, that learns multiple clusterings simultaneously, and is trained fully unsupervised and end-to-end. MFCVAE uses a progressively-trained ladder architecture which leads to highly stable performance. We provide novel theoretical results for optimising the ELBO analytically with respect to the categorical variational posterior distribution, correcting earlier influential theoretical work. On image benchmarks, we demonstrate that our approach separates out and clusters over different aspects of the data in a disentangled manner. We also show other advantages of our model: the compositionality of its latent space and that it provides controlled generation of samples.

[Synthetic Design: An Optimization Approach to Experimental Design with Synthetic Controls](#)

- Nick Doudchenko, Khashayar Khosravi, Jean Pouget-Abadie, Sébastien Lahaie, Miles Lubin, Vahab Mirrokni, Jann Spiess, guido imbens
- abstract@[open-review](#): We investigate the optimal design of experimental studies that have pre-treatment outcome data available. The average treatment effect is estimated as the difference between the weighted average outcomes of the treated and control units. A number of commonly used approaches fit this formulation, including the difference-in-means estimator and a variety of synthetic-control techniques. We propose several methods for choosing the set of treated units in conjunction with the weights. Observing the NP-hardness of the problem, we introduce a mixed-integer programming formulation which selects both the treatment and control sets and unit weightings. We prove that these proposed approaches lead to qualitatively different experimental units being selected for treatment. We use simulations based on publicly available data from the US Bureau of Labor Statistics that show improvements in terms of mean squared error and statistical power when compared to simple and commonly used alternatives such as randomized trials.

[Ranking Policy Decisions](#)

- Hadrien Pouget, Hana Chockler, Youcheng Sun, Daniel Kroening
- abstract@[open-review](#): Policies trained via Reinforcement Learning (RL) without human intervention are often needlessly complex, making them difficult to analyse and interpret. In a run with $\$n\$$ time steps, a policy will make $\$n\$$ decisions on actions to take; we conjecture that only a small subset of these decisions delivers value over selecting a simple default action. Given a trained policy, we propose a novel black-box method based on statistical fault localisation that ranks the states of the environment according to the importance of decisions made in those states. We argue that among other things, the ranked list of states can help explain and understand the policy. As the ranking method is statistical, a direct evaluation of its quality is hard. As a proxy for quality, we use the ranking to create new, simpler policies from the original ones by pruning decisions identified as unimportant (that is, replacing them by default actions) and measuring the impact on performance. Our experimental results on a diverse set of standard benchmarks demonstrate that pruned policies can perform on a level comparable to the original policies. We show that naive approaches for ranking policies, e.g. ranking based on the frequency of visiting a state, do not result in high-performing pruned policies. To the best of our knowledge, there are no similar techniques for ranking RL policies' decisions.

[Searching the Search Space of Vision Transformer](#)

- Minghao Chen, Kan Wu, Bolin Ni, Houwen Peng, Bei Liu, Jianlong Fu, Hongyang Chao, Haibin Ling
- abstract@[open-review](#): Vision Transformer has shown great visual representation power in substantial vision tasks such as recognition and detection, and thus been attracting fast-growing efforts on manually designing more effective architectures. In this paper, we propose to use neural architecture search to automate this process, by searching not only the architecture but also the search space. The central idea is to gradually evolve different search dimensions guided by their E-T Error computed using a weight-sharing supernet. Moreover, we provide design guidelines of general vision transformers with extensive analysis according to the space searching process, which could promote the understanding of vision transformer. Remarkably, the searched models, named S3 (short for Searching the Search Space), from the searched space achieve superior performance to recently proposed models, such as Swin, DeiT and ViT, when evaluated on ImageNet. The effectiveness of S3 is also illustrated on object detection, semantic segmentation and visual question answering, demonstrating its generality to downstream vision and vision-language tasks. Code and models will be available at <https://github.com/microsoft/Cream>.

[Relative stability toward diffeomorphisms indicates performance in deep nets](#)

- Leonardo Petrini, Alessandro Favero, Mario Geiger, Matthieu Wyart
- abstract@[open-review](#): Understanding why deep nets can classify data in large dimensions remains a challenge. It has been proposed that they do so by becoming stable to diffeomorphisms, yet existing empirical measurements support that it is often not the case. We revisit this question by defining a maximum-entropy distribution on diffeomorphisms, that allows to study typical diffeomorphisms of a given norm. We confirm that stability toward

diffeomorphisms does not strongly correlate to performance on benchmark data sets of images. By contrast, we find that the stability toward diffeomorphisms relative to that of generic transformations R_f correlates remarkably with the test error ϵ_t . It is of order unity at initialization but decreases by several decades during training for state-of-the-art architectures. For CIFAR10 and 15 known architectures, we find $\epsilon_t \approx 0.2\sqrt{R_f}$, suggesting that obtaining a small R_f is important to achieve good performance. We study how R_f depends on the size of the training set and compare it to a simple model of invariant learning.

[Raw Nav-merge Seismic Data to Subsurface Properties with MLP based Multi-Modal Information Unscrambler](#)

- Aditya Desai, Zhaozhuo Xu, Menal Gupta, Anu Chandran, Antoine Vial-Aussavy, Anshumali Shrivastava
- abstract@[open-review](#): Traditional seismic inversion (SI) maps the hundreds of terabytes of raw-field data to subsurface properties in gigabytes. This inversion process is expensive, requiring over a year of human and computational effort. Recently, data-driven approaches equipped with Deep learning (DL) are envisioned to improve SI efficiency. However, these improvements are restricted to data with highly reduced scale and complexity. To extend these approaches to real-scale seismic data, researchers need to process raw nav-merge seismic data into an image and perform convolution. We argue that this convolution-based way of SI is not only computationally expensive but also conceptually problematic. Seismic data is not naturally an image and need not be processed as images. In this work, we go beyond convolution and propose a novel SI method. We solve the scalability of SI by proposing a new auxiliary learning paradigm for SI (Aux-SI). This paradigm breaks the SI into local inversion tasks, which predicts each small chunk of subsurface properties using surrounding seismic data. Aux-SI combines these local predictions to obtain the entire subsurface model. However, even this local inversion is still challenging due to: (1) high-dimensional, spatially irregular multi-modal seismic data, (2) there is no concrete spatial mapping (or alignment) between subsurface properties and raw data. To handle these challenges, we propose an all-MLP architecture, Multi-Modal Information Unscrambler (MMI-Unscrambler), that unscrambles seismic information by ingesting all available multi-modal data. The experiment shows that MMI-Unscrambler outperforms both SOTA U-Net and Transformer models on simulation data. We also scale MMI-Unscrambler to raw-field nav-merge data on Gulf-of-Mexico to obtain a geologically sound velocity model with an SSIM score of 0.8. To the best of our knowledge, this is the first successful demonstration of the DL approach on SI for real, large-scale, and complicated raw field data.

[Inverse Problems Leveraging Pre-trained Contrastive Representations](#)

- Sriram Ravula, Georgios Smyrnis, Matt Jordan, Alexandros G. Dimakis
- abstract@[open-review](#): We study a new family of inverse problems for recovering representations of corrupted data. We assume access to a pre-trained representation learning network $R(x)$ that operates on clean images, like CLIP. The problem is to recover the representation of an image $R(x)$, if we are only given a corrupted version $A(x)$, for some known forward operator A . We propose a supervised inversion method that uses a contrastive objective to obtain excellent representations for highly corrupted images. Using a linear probe on our robust representations, we achieve a higher accuracy than end-to-end supervised baselines when classifying images with various types of distortions, including blurring, additive noise, and random pixel masking. We evaluate on a subset of ImageNet and observe that our method is robust to varying levels of distortion. Our method outperforms end-to-end baselines even with a fraction of the labeled data in a wide range of forward operators.

[The Unbalanced Gromov Wasserstein Distance: Conic Formulation and Relaxation](#)

- Thibault Sejourne, Francois-Xavier Vialard, Gabriel Peyré
- abstract@[open-review](#): Comparing metric measure spaces (i.e. a metric space endowed with a probability distribution) is at the heart of many machine learning problems. The most popular distance between such metric measure spaces is the Gromov-Wasserstein (GW) distance, which is the solution of a quadratic assignment problem. The GW distance is however limited to the comparison of metric measure spaces endowed with a probability distribution. To alleviate this issue, we introduce two Unbalanced Gromov-Wasserstein formulations: a distance and a more tractable upper-bounding relaxation. They both allow the comparison of metric spaces equipped with arbitrary positive measures up to isometries. The first formulation is a positive and definite divergence based on a relaxation of the mass conservation constraint using a novel type of quadratically-homogeneous divergence. This divergence works hand in hand with the entropic regularization approach which is popular to solve large scale optimal transport problems. We show that the underlying non-convex optimization problem can be efficiently tackled using a highly parallelizable and GPU-friendly iterative scheme. The second formulation is a distance between mm-spaces up to isometries based on a conic lifting. Lastly, we provide numerical experiments on synthetic and domain adaptation data with a Positive-Unlabeled learning task to highlight the salient features of the unbalanced divergence and its potential applications in ML.

[Diffusion Models Beat GANs on Image Synthesis](#)

- Prafulla Dhariwal, Alexander Nichol
- abstract@[open-review](#): We show that diffusion models can achieve image sample quality superior to the current state-of-the-art generative models. We achieve this on unconditional image synthesis by finding a better architecture through a series of ablations. For conditional image synthesis, we further improve sample quality with classifier guidance: a simple, compute-efficient method for trading off diversity for fidelity using gradients from a classifier. We achieve an FID of 2.97 on ImageNet 128\$\times\$128, 4.59 on ImageNet 256\$\times\$256, and 7.72 on ImageNet 512\$\times\$512, and we match BigGAN-deep even with as few as 25 forward passes per sample, all while maintaining better coverage of the distribution. Finally, we find that classifier guidance combines well with upsampling diffusion models, further improving FID to 3.94 on ImageNet 256\$\times\$256 and 3.85 on ImageNet 512\$\times\$512.

[Learning MDPs from Features: Predict-Then-Optimize for Sequential Decision Making by Reinforcement Learning](#)

- Kai Wang, Sanket Shah, Haipeng Chen, Andrew Perrault, Finale Doshi-Velez, Milind Tambe
- abstract@[open-review](#): In the predict-then-optimize framework, the objective is to train a predictive model, mapping from environment features to parameters of an optimization problem, which maximizes decision quality when the optimization is subsequently solved. Recent work on decision-focused learning shows that embedding the optimization problem in the training pipeline can improve decision quality and help generalize better to unseen tasks compared to relying on an intermediate loss function for evaluating prediction quality. We study the predict-then-optimize framework in the context of sequential decision problems (formulated as MDPs) that are solved via reinforcement learning. In particular, we are given environment features and a set of trajectories from training MDPs, which we use to train a predictive model that generalizes to unseen test MDPs without trajectories. Two significant computational challenges arise in applying decision-focused learning to MDPs: (i) large state and action spaces make it infeasible for existing techniques to differentiate through MDP problems, and (ii) the high-dimensional policy space, as parameterized by a neural network, makes differentiating through a policy expensive. We resolve the first challenge by sampling provably unbiased derivatives to approximate and differentiate through optimality conditions, and the second challenge by using a low-rank approximation to the high-dimensional sample-based derivatives. We implement both Bellman-based and policy gradient-based decision-focused learning on three different MDP problems with missing parameters, and show that decision-focused learning performs better in generalization to unseen tasks.

[A Closer Look at the Worst-case Behavior of Multi-armed Bandit Algorithms](#)

- Anand Kalvit, Assaf Zeevi
- abstract@[open-review](#): One of the key drivers of complexity in the classical (stochastic) multi-armed bandit (MAB) problem is the difference between mean rewards in the top two arms, also known as the instance gap. The celebrated Upper Confidence Bound (UCB) policy is among the simplest optimism-based MAB algorithms that naturally adapts to this gap: for a horizon of play n , it achieves optimal $O(\log n)$ regret in instances with "large"

gaps, and a near-optimal $O(\sqrt{n \log n})$ minimax regret when the gap can be arbitrarily "small." This paper provides new results on the arm-sampling behavior of UCB, leading to several important insights. Among these, it is shown that arm-sampling rates under UCB are asymptotically deterministic, regardless of the problem complexity. This discovery facilitates new sharp asymptotics and a novel alternative proof for the $O(\sqrt{n \log n})$ minimax regret of UCB. Furthermore, the paper also provides the first complete process-level characterization of the MAB problem in the conventional diffusion scaling. Among other things, the "small" gap worst-case lens adopted in this paper also reveals profound distinctions between the behavior of UCB and Thompson Sampling, such as an "incomplete learning" phenomenon characteristic of the latter.

[SAPE: Spatially-Adaptive Progressive Encoding for Neural Optimization](#)

- Amir Hertz, Or Perel, Raja Giryes, Olga Sorkine-hornung, Daniel Cohen-or
- abstract@[open-review](#): Multilayer-perceptrons (MLP) are known to struggle learning functions of high-frequencies, and in particular, instances of wide frequency bands. We present a progressive mapping scheme for input signals of MLP networks, enabling them to better fit a wide range of frequencies without sacrificing training stability or requiring any domain specific preprocessing. We introduce Spatially Adaptive Progressive Encoding (SAPE) layers, which gradually unmask signal components with increasing frequencies as a function of time and space. The progressive exposure of frequencies is monitored by a feedback loop throughout the neural optimization process, allowing changes to propagate at different rates among local spatial portions of the signal space. We demonstrate the advantage of our method on variety of domains and applications: regression of low dimensional signals and images, representation learning of occupancy networks, and a geometric task of mesh transfer between 3D shapes.

[A Biased Graph Neural Network Sampler with Near-Optimal Regret](#)

- Qingru Zhang, David Wipf, Quan Gan, Le Song
- abstract@[open-review](#): Graph neural networks (GNN) have recently emerged as a vehicle for applying deep network architectures to graph and relational data. However, given the increasing size of industrial datasets, in many practical situations, the message passing computations required for sharing information across GNN layers are no longer scalable. Although various sampling methods have been introduced to approximate full-graph training within a tractable budget, there remain unresolved complications such as high variances and limited theoretical guarantees. To address these issues, we build upon existing work and treat GNN neighbor sampling as a multi-armed bandit problem but with a newly-designed reward function that introduces some degree of bias designed to reduce variance and avoid unstable, possibly-unbounded pay outs. And unlike prior bandit-GNN use cases, the resulting policy leads to near-optimal regret while accounting for the GNN training dynamics introduced by SGD. From a practical standpoint, this translates into lower variance estimates and competitive or superior test accuracy across several benchmarks.

[Equilibrium Refinement for the Age of Machines: The One-Sided Quasi-Perfect Equilibrium](#)

- Gabriele Farina, Tuomas Sandholm
- abstract@[open-review](#): In two-player zero-sum extensive-form games, Nash equilibrium prescribes optimal strategies against perfectly rational opponents. However, it does not guarantee rational play in parts of the game tree that can only be reached by the players making mistakes. This can be problematic when operationalizing equilibria in the real world among imperfect players. Trembling-hand refinements are a sound remedy to this issue, and are subsets of Nash equilibria that are designed to handle the possibility that any of the players may make mistakes. In this paper, we initiate the study of equilibrium refinements for settings where one of the players is perfectly rational (the ``machine'') and the other may make mistakes. As we show, this endeavor has many pitfalls: many intuitively appealing approaches to refinement fail in various ways. On the positive side, we introduce a modification of the classical quasi-perfect equilibrium (QPE) refinement, which we call the one-sided quasi-perfect equilibrium. Unlike QPE, one-sided QPE only accounts for mistakes from one player and assumes that no mistakes will be made by the machine. We present experiments on standard benchmark games and an endgame from the famous man-machine match where the AI Libratus was the first to beat top human specialist professionals in heads-up no-limit Texas hold'em poker. We show that one-sided QPE can be computed more efficiently than all known prior refinements, paving the way to wider adoption of Nash equilibrium refinements in settings with perfectly rational machines (or humans perfectly actuating machine-generated strategies) that interact with players prone to mistakes. We also show that one-sided QPE tends to play better than a Nash equilibrium strategy against imperfect opponents.

[Interpreting Representation Quality of DNNs for 3D Point Cloud Processing](#)

- Wen Shen, Qihan Ren, Dongrui Liu, Quanshi Zhang
- abstract@[open-review](#): In this paper, we evaluate the quality of knowledge representations encoded in deep neural networks (DNNs) for 3D point cloud processing. We propose a method to disentangle the overall model vulnerability into the sensitivity to the rotation, the translation, the scale, and local 3D structures. Besides, we also propose metrics to evaluate the spatial smoothness of encoding 3D structures, and the representation complexity of the DNN. Based on such analysis, experiments expose representation problems with classic DNNs, and explain the utility of the adversarial training. The code will be released when this paper is accepted.

[How Fine-Tuning Allows for Effective Meta-Learning](#)

- Kurtland Chua, Qi Lei, Jason D. Lee
- abstract@[open-review](#): Representation learning has served as a key tool for meta-learning, enabling rapid learning of new tasks. Recent works like MAML learn task-specific representations by finding an initial representation requiring minimal per-task adaptation (i.e. a fine-tuning-based objective). We present a theoretical framework for analyzing a MAML-like algorithm, assuming all available tasks require approximately the same representation. We then provide risk bounds on predictors found by fine-tuning via gradient descent, demonstrating that the method provably leverages the shared structure. We illustrate these bounds in the logistic regression and neural network settings. In contrast, we establish settings where learning one representation for all tasks (i.e. using a "frozen representation" objective) fails. Notably, any such algorithm cannot outperform directly learning the target task with no other information, in the worst case. This separation underscores the benefit of fine-tuning-based over "frozen representation" objectives in few-shot learning.

[Cooperative Stochastic Bandits with Asynchronous Agents and Constrained Feedback](#)

- Lin Yang, Yu-Zhen Janice Chen, Stephen Pasteris, Mohammad Hajiesmaili, John C. S. Lui, Don Towsley
- abstract@[open-review](#): This paper studies a cooperative multi-armed bandit problem with \$M\$ agents cooperating together to solve the same instance of a \$K\$-armed stochastic bandit problem with the goal of maximizing the cumulative reward of agents. The agents are heterogeneous in (i) their limited access to a local subset of arms; and (ii) their decision-making rounds, i.e., agents are asynchronous with different decision-making gaps. The goal is to find the global optimal arm and agents are able to pull any arm, however, they observe the reward only when the selected arm is local. The challenge is a tradeoff for agents between pulling a local arm with the possibility of observing the feedback, or relying on the observations of other agents that might occur at different rates. Naive extensions of traditional algorithms lead to an arbitrarily poor regret as a function of aggregate action frequency of any \$suboptimal\$ arm located in slow agents. We resolve this issue by proposing a novel two-stage learning algorithm, called \$\text{CO-LCB}\$, whose regret is a function of aggregate action frequency of agents containing the \$\text{optimal}\$ arm. We also show that the regret of \$\text{CO-LCB}\$ matches the regret lower bound up to a small factor.

[Multiple Descent: Design Your Own Generalization Curve](#)

- Lin Chen, Yifei Min, Mikhail Belkin, Amin Karbasi
- abstract@[open-review](#): This paper explores the generalization loss of linear regression in variably parameterized families of models, both under-parameterized and over-parameterized. We show that the generalization curve can have an arbitrary number of peaks, and moreover, the locations of those peaks can be explicitly controlled. Our results highlight the fact that both the classical U-shaped generalization curve and the recently observed double descent curve are not intrinsic properties of the model family. Instead, their emergence is due to the interaction between the properties of the data and the inductive biases of learning algorithms.

[On Empirical Risk Minimization with Dependent and Heavy-Tailed Data](#)

- Abhishek Roy, Krishnakumar Balasubramanian, Murat A. Erdogdu
- abstract@[open-review](#): In this work, we establish risk bounds for Empirical Risk Minimization (ERM) with both dependent and heavy-tailed data-generating processes. We do so by extending the seminal works~\cite{pmlr-v35-mendelson14, mendelson2018learning} on the analysis of ERM with heavy-tailed but independent and identically distributed observations, to the strictly stationary exponentially β -mixing case. We allow for the interaction between the noise and inputs to be even polynomially heavy-tailed, which covers a significantly large class of heavy-tailed models beyond what is analyzed in the learning theory literature. We illustrate our theoretical results by obtaining rates of convergence for high-dimensional linear regression with dependent and heavy-tailed data.

[Gone Fishing: Neural Active Learning with Fisher Embeddings](#)

- Jordan Ash, Surbhi Goel, Akshay Krishnamurthy, Sham Kakade
- abstract@[open-review](#): There is an increasing need for effective active learning algorithms that are compatible with deep neural networks. This paper motivates and revisits a classic, Fisher-based active selection objective, and proposes BAIT, a practical, tractable, and high-performing algorithm that makes it viable for use with neural models. BAIT draws inspiration from the theoretical analysis of maximum likelihood estimators (MLE) for parametric models. It selects batches of samples by optimizing a bound on the MLE error in terms of the Fisher information, which we show can be implemented efficiently at scale by exploiting linear-algebraic structure especially amenable to execution on modern hardware. Our experiments demonstrate that BAIT outperforms the previous state of the art on both classification and regression problems, and is flexible enough to be used with a variety of model architectures.

[On Riemannian Optimization over Positive Definite Matrices with the Bures-Wasserstein Geometry](#)

- Andi Han, Bamdev Mishra, Pratik Kumar Jawanpuria, Junbin Gao
- abstract@[open-review](#): In this paper, we comparatively analyze the Bures-Wasserstein (BW) geometry with the popular Affine-Invariant (AI) geometry for Riemannian optimization on the symmetric positive definite (SPD) matrix manifold. Our study begins with an observation that the BW metric has a linear dependence on SPD matrices in contrast to the quadratic dependence of the AI metric. We build on this to show that the BW metric is a more suitable and robust choice for several Riemannian optimization problems over ill-conditioned SPD matrices. We show that the BW geometry has a non-negative curvature, which further improves convergence rates of algorithms over the non-positively curved AI geometry. Finally, we verify that several popular cost functions, which are known to be geodesic convex under the AI geometry, are also geodesic convex under the BW geometry. Extensive experiments on various applications support our findings.

[Refining Language Models with Compositional Explanations](#)

- Huihan Yao, Ying Chen, Qinyuan Ye, Xisen Jin, Xiang Ren
- abstract@[open-review](#): Pre-trained language models have been successful on text classification tasks, but are prone to learning spurious correlations from biased datasets, and are thus vulnerable when making inferences in a new domain. Prior work reveals such spurious patterns via post-hoc explanation algorithms which compute the importance of input features. Further, the model is regularized to align the importance scores with human knowledge, so that the unintended model behaviors are eliminated. However, such a regularization technique lacks flexibility and coverage, since only importance scores towards a pre-defined list of features are adjusted, while more complex human knowledge such as feature interaction and pattern generalization can hardly be incorporated. In this work, we propose to refine a learned language model for a target domain by collecting human-provided compositional explanations regarding observed biases. By parsing these explanations into executable logic rules, the human-specified refinement advice from a small set of explanations can be generalized to more training examples. We additionally introduce a regularization term allowing adjustments for both importance and interaction of features to better rectify model behavior. We demonstrate the effectiveness of the proposed approach on two text classification tasks by showing improved performance in target domain as well as improved model fairness after refinement.

[Going Beyond Linear RL: Sample Efficient Neural Function Approximation](#)

- Baihe Huang, Kaixuan Huang, Sham Kakade, Jason D. Lee, Qi Lei, Runzhe Wang, Jiaqi Yang
- abstract@[open-review](#): Deep Reinforcement Learning (RL) powered by neural net approximation of the Q function has had enormous empirical success. While the theory of RL has traditionally focused on linear function approximation (or eluder dimension) approaches, little is known about nonlinear RL with neural net approximations of the Q functions. This is the focus of this work, where we study function approximation with two-layer neural networks (considering both ReLU and polynomial activation functions). Our first result is a computationally and statistically efficient algorithm in the generative model setting under completeness for two-layer neural networks. Our second result considers this setting but under only realizability of the neural net function class. Here, assuming deterministic dynamics, the sample complexity scales linearly in the algebraic dimension. In all cases, our results significantly improve upon what can be attained with linear (or eluder dimension) methods.

[Scalable Neural Data Server: A Data Recommender for Transfer Learning](#)

- Tianshi Cao, Sasha (Alexandre) Doubov, David Acuna, Sanja Fidler
- abstract@[open-review](#): Absence of large-scale labeled data in the practitioner's target domain can be a bottleneck to applying machine learning algorithms in practice. Transfer learning is a popular strategy for leveraging additional data to improve the downstream performance, but finding the most relevant data to transfer from can be challenging. Neural Data Server (NDS), a search engine that recommends relevant data for a given downstream task, has been previously proposed to address this problem (Yan et al., 2020). NDS uses a mixture of experts trained on data sources to estimate similarity between each source and the downstream task. Thus, the computational cost to each user grows with the number of sources and requires an expensive training step for each data provider. To address these issues, we propose Scalable Neural Data Server (SNDS), a large-scale search engine that can theoretically index thousands of datasets to serve relevant ML data to end users. SNDS trains the mixture of experts on intermediary datasets during initialization, and represents both data sources and downstream tasks by their proximity to the intermediary datasets. As such, computational cost incurred by users of SNDS remains fixed as new datasets are added to the server, without pre-training for the data providers. We validate SNDS on a plethora of real world tasks and find that data recommended by SNDS improves downstream task performance over baselines. We also demonstrate the scalability of our system by demonstrating its ability to select relevant data for transfer outside of the natural image setting.

[What can linearized neural networks actually say about generalization?](#)

- Guillermo Ortiz-Jimenez, Seyed-Mohsen Moosavi-Dezfooli, Pascal Frossard
- abstract@[open-review](#): For certain infinitely-wide neural networks, the neural tangent kernel (NTK) theory fully characterizes generalization, but for the networks used in practice, the empirical NTK only provides a rough first-order approximation. Still, a growing body of work keeps leveraging this approximation to successfully analyze important deep learning phenomena and design algorithms for new applications. In our work, we provide strong empirical evidence to determine the practical validity of such approximation by conducting a systematic comparison of the behavior of different neural networks and their linear approximations on different tasks. We show that the linear approximations can indeed rank the learning complexity of certain tasks for neural networks, even when they achieve very different performances. However, in contrast to what was previously reported, we discover that neural networks do not always perform better than their kernel approximations, and reveal that the performance gap heavily depends on architecture, dataset size and training task. We discover that networks overfit to these tasks mostly due to the evolution of their kernel during training, thus, revealing a new type of implicit bias.

[CATs: Cost Aggregation Transformers for Visual Correspondence](#)

- Seokju Cho, Sunghwan Hong, Sangryul Jeon, Yunsung Lee, Kwanghoon Sohn, Seungryong Kim
- abstract@[open-review](#): We propose a novel cost aggregation network, called Cost Aggregation Transformers (CATs), to find dense correspondences between semantically similar images with additional challenges posed by large intra-class appearance and geometric variations. Cost aggregation is a highly important process in matching tasks, which the matching accuracy depends on the quality of its output. Compared to hand-crafted or CNN-based methods addressing the cost aggregation, in that either lacks robustness to severe deformations or inherit the limitation of CNNs that fail to discriminate incorrect matches due to limited receptive fields, CATs explore global consensus among initial correlation map with the help of some architectural designs that allow us to fully leverage self-attention mechanism. Specifically, we include appearance affinity modeling to aid the cost aggregation process in order to disambiguate the noisy initial correlation maps and propose multi-level aggregation to efficiently capture different semantics from hierarchical feature representations. We then combine with swapping self-attention technique and residual connections not only to enforce consistent matching, but also to ease the learning process, which we find that these result in an apparent performance boost. We conduct experiments to demonstrate the effectiveness of the proposed model over the latest methods and provide extensive ablation studies. Code and trained models are available at <https://sunghwanhong.github.io/CATs/>.

[Asynchronous Stochastic Optimization Robust to Arbitrary Delays](#)

- Alon Cohen, Amit Daniely, Yoel Drori, Tomer Koren, Mariano Schain
- abstract@[open-review](#): We consider the problem of stochastic optimization with delayed gradients in which, at each time step t , the algorithm makes an update using a stale stochastic gradient from step $t - d_t$ for some arbitrary delay d_t . This setting abstracts asynchronous distributed optimization where a central server receives gradient updates computed by worker machines. These machines can experience computation and communication loads that might vary significantly over time. In the general non-convex smooth optimization setting, we give a simple and efficient algorithm that requires $O(\sigma^2/\epsilon^4 + \tau/\epsilon^2)$ steps for finding an ϵ -stationary point x . Here, τ is the average delay $\frac{1}{T} \sum_{t=1}^T d_t$ and σ^2 is the variance of the stochastic gradients. This improves over previous work, which showed that stochastic gradient descent achieves the same rate but with respect to the maximal delay $\max_t d_t$, that can be significantly larger than the average delay especially in heterogeneous distributed systems. Our experiments demonstrate the efficacy and robustness of our algorithm in cases where the delay distribution is skewed or heavy-tailed.

[Consistent Non-Parametric Methods for Maximizing Robustness](#)

- Robi Bhattacharjee, Kamalika Chaudhuri
- abstract@[open-review](#): Learning classifiers that are robust to adversarial examples has received a great deal of recent attention. A major drawback of the standard robust learning framework is the imposition of an artificial robustness radius r that applies to all inputs, and ignores the fact that data may be highly heterogeneous. In particular, it is plausible that robustness regions should be larger in some regions of data, and smaller in other. In this paper, we address this limitation by proposing a new limit classifier, called the neighborhood optimal classifier, that extends the Bayes optimal classifier outside its support by using the label of the closest in-support point. We then argue that this classifier maximizes the size of its robustness regions subject to the constraint of having accuracy equal to the Bayes optimal. We then present sufficient conditions under which general non-parametric methods that can be represented as weight functions converge towards this limit object, and show that both nearest neighbors and kernel classifiers (under certain assumptions) suffice.

[Generalizable Multi-linear Attention Network](#)

- Tao Jin, Zhou Zhao
- abstract@[open-review](#): The majority of existing multimodal sequential learning methods focus on how to obtain effective representations and ignore the importance of multimodal fusion. Bilinear attention network (BAN) is a commonly used fusion method, which leverages tensor operations to associate the features of different modalities. However, BAN has a poor compatibility for more modalities, since the computational complexity of the attention map increases exponentially with the number of modalities. Based on this concern, we propose a new method called generalizable multi-linear attention network (MAN), which can associate as many modalities as possible in linear complexity with hierarchical approximation decomposition (HAD). Besides, considering the fact that softmax attention kernels cannot be decomposed as linear operation directly, we adopt the addition random features (ARF) mechanism to approximate the non-linear softmax functions with enough theoretical analysis. We conduct extensive experiments on four datasets of three tasks (multimodal sentiment analysis, multimodal speaker traits recognition, and video retrieval), the experimental results show that MAN could achieve competitive results compared with the state-of-the-art methods, showcasing the effectiveness of the approximation decomposition and addition random features mechanism.

[Labeling Trick: A Theory of Using Graph Neural Networks for Multi-Node Representation Learning](#)

- Muhan Zhang, Pan Li, Yinglong Xia, Kai Wang, Long Jin
- abstract@[open-review](#): In this paper, we provide a theory of using graph neural networks (GNNs) for multi-node representation learning (where we are interested in learning a representation for a set of more than one node, such as link). We know that GNN is designed to learn single-node representations. When we want to learn a node set representation involving multiple nodes, a common practice in previous works is to directly aggregate the single-node representations obtained by a GNN into a joint node set representation. In this paper, we show a fundamental constraint of such an approach, namely the inability to capture the dependence between nodes in the node set, and argue that directly aggregating individual node representations does not lead to an effective joint representation for multiple nodes. Then, we notice that a few previous successful works for multi-node representation learning, including SEAL, Distance Encoding, and ID-GNN, all used node labeling. These methods first label nodes in the graph according to their relationships with the target node set before applying a GNN. Then, the node representations obtained in the labeled graph are aggregated into a node set representation. By investigating their inner mechanisms, we unify these node labeling techniques into a single and most general form---labeling trick. We prove that with labeling trick a sufficiently expressive GNN learns the most expressive node set representations, thus in principle solves any joint learning tasks over node sets. Experiments on one important two-node representation learning task, link prediction, verified our theory. Our work explains the superior performance of previous node-labeling-based methods, and establishes a theoretical foundation of using GNNs for multi-node representation learning.

[SUPER-ADAM: Faster and Universal Framework of Adaptive Gradients](#)

- Feihu Huang, Junyi Li, Heng Huang
- abstract@[open-review](#): Adaptive gradient methods have shown excellent performances for solving many machine learning problems. Although multiple adaptive gradient methods were recently studied, they mainly focus on either empirical or theoretical aspects and also only work for specific problems by using some specific adaptive learning rates. Thus, it is desired to design a universal framework for practical algorithms of adaptive gradients with theoretical guarantee to solve general problems. To fill this gap, we propose a faster and universal framework of adaptive gradients (i.e., SUPER-ADAM) by introducing a universal adaptive matrix that includes most existing adaptive gradient forms. Moreover, our framework can flexibly integrate the momentum and variance reduced techniques. In particular, our novel framework provides the convergence analysis support for adaptive gradient methods under the nonconvex setting. In theoretical analysis, we prove that our SUPER-ADAM algorithm can achieve the best known gradient (i.e., stochastic first-order oracle (SFO)) complexity of $\tilde{O}(\epsilon^{-3})$ for finding an ϵ -stationary point of nonconvex optimization, which matches the lower bound for stochastic smooth nonconvex optimization. In numerical experiments, we employ various deep learning tasks to validate that our algorithm consistently outperforms the existing adaptive algorithms. Code is available at <https://github.com/LIJUNYI95/SuperAdam>

[General Nonlinearities in SO\(2\)-Equivariant CNNs](#)

- Daniel Franzen, Michael Wand
- abstract@[open-review](#): Invariance under symmetry is an important problem in machine learning. Our paper looks specifically at equivariant neural networks where transformations of inputs yield homomorphic transformations of outputs. Here, steerable CNNs have emerged as the standard solution. An inherent problem of steerable representations is that general nonlinear layers break equivariance, thus restricting architectural choices. Our paper applies harmonic distortion analysis to illuminate the effect of nonlinearities on Fourier representations of $SO(2)$. We develop a novel FFT-based algorithm for computing representations of non-linearly transformed activations while maintaining band-limitation. It yields exact equivariance for polynomial (approximations of) nonlinearities, as well as approximate solutions with tunable accuracy for general functions. We apply the approach to build a fully $E(3)$ -equivariant network for sampled 3D surface data. In experiments with 2D and 3D data, we obtain results that compare favorably to the state-of-the-art in terms of accuracy while permitting continuous symmetry and exact equivariance.

[Denoising Normalizing Flow](#)

- Christian Horvat, Jean-Pascal Pfister
- abstract@[open-review](#): Normalizing flows (NF) are expressive as well as tractable density estimation methods whenever the support of the density is diffeomorphic to the entire data-space. However, real-world data sets typically live on (or very close to) low-dimensional manifolds thereby challenging the applicability of standard NF on real-world problems. Here we propose a novel method - called Denoising Normalizing Flow (DNF) - that estimates the density on the low-dimensional manifold while learning the manifold as well. The DNF works in 3 steps. First, it inflates the manifold - making it diffeomorphic to the entire data-space. Secondly, it learns an NF on the inflated manifold and finally it learns a denoising mapping - similarly to denoising autoencoders. The DNF relies on a single cost function and does not require to alternate between a density estimation phase and a manifold learning phase - as it is the case with other recent methods. Furthermore, we show that the DNF can learn meaningful low-dimensional representations from naturalistic images as well as generate high-quality samples.

[Attention over Learned Object Embeddings Enables Complex Visual Reasoning](#)

- David Ding, Felix Hill, Adam Santoro, Malcolm Reynolds, Matt Botvinick
- abstract@[open-review](#): Neural networks have achieved success in a wide array of perceptual tasks but often fail at tasks involving both perception and higher-level reasoning. On these more challenging tasks, bespoke approaches (such as modular symbolic components, independent dynamics models or semantic parsers) targeted towards that specific type of task have typically performed better. The downside to these targeted approaches, however, is that they can be more brittle than general-purpose neural networks, requiring significant modification or even redesign according to the particular task at hand. Here, we propose a more general neural-network-based approach to dynamic visual reasoning problems that obtains state-of-the-art performance on three different domains, in each case outperforming bespoke modular approaches tailored specifically to the task. Our method relies on learned object-centric representations, self-attention and self-supervised dynamics learning, and all three elements together are required for strong performance to emerge. The success of this combination suggests that there may be no need to trade off flexibility for performance on problems involving spatio-temporal or causal-style reasoning. With the right soft biases and learning objectives in a neural network we may be able to attain the best of both worlds.

[Differentially Private Federated Bayesian Optimization with Distributed Exploration](#)

- Zhongxiang Dai, Bryan Kian Hsiang Low, Patrick Jaillet
- abstract@[open-review](#): Bayesian optimization (BO) has recently been extended to the federated learning (FL) setting by the federated Thompson sampling (FTS) algorithm, which has promising applications such as federated hyperparameter tuning. However, FTS is not equipped with a rigorous privacy guarantee which is an important consideration in FL. Recent works have incorporated differential privacy (DP) into the training of deep neural networks through a general framework for adding DP to iterative algorithms. Following this general DP framework, our work here integrates DP into FTS to preserve user-level privacy. We also leverage the ability of this general DP framework to handle different parameter vectors, as well as the technique of local modeling for BO, to further improve the utility of our algorithm through distributed exploration (DE). The resulting differentially private FTS with DE (DP-FTS-DE) algorithm is endowed with theoretical guarantees for both the privacy and utility and is amenable to interesting theoretical insights about the privacy-utility trade-off. We also use real-world experiments to show that DP-FTS-DE achieves high utility (competitive performance) with a strong privacy guarantee (small privacy loss) and induces a trade-off between privacy and utility.

[Differentiable Learning Under Triage](#)

- Nastaran Okati, Abir De, Manuel Rodriguez
- abstract@[open-review](#): Multiple lines of evidence suggest that predictive models may benefit from algorithmic triage. Under algorithmic triage, a predictive model does not predict all instances but instead defers some of them to human experts. However, the interplay between the prediction accuracy of the model and the human experts under algorithmic triage is not well understood. In this work, we start by formally characterizing under which circumstances a predictive model may benefit from algorithmic triage. In doing so, we also demonstrate that models trained for full automation may be suboptimal under triage. Then, given any model and desired level of triage, we show that the optimal triage policy is a deterministic threshold rule in which triage decisions are derived deterministically by thresholding the difference between the model and human errors on a per-instance level. Building upon these results, we introduce a practical gradient-based algorithm that is guaranteed to find a sequence of predictive models and triage policies of increasing performance. Experiments on a wide variety of supervised learning tasks using synthetic and real data from two important applications---content moderation and scientific discovery---illustrate our theoretical results and show that the models and triage policies provided by our gradient-based algorithm outperform those provided by several competitive baselines.

[ROI Maximization in Stochastic Online Decision-Making](#)

- Nicolò Cesa-Bianchi, Tom Cesari, Yishay Mansour, Vianney Perchet
- abstract@[open-review](#): We introduce a novel theoretical framework for Return On Investment (ROI) maximization in repeated decision-making. Our setting is motivated by the use case of companies that regularly receive proposals for technological innovations and want to quickly decide whether they are worth implementing. We design an algorithm for learning ROI-maximizing decision-making policies over a sequence of innovation proposals. Our

algorithm provably converges to an optimal policy in class Π at a rate of order $\min\{\frac{1}{(N\Delta)^2}, N^{-1/3}\}$, where N is the number of innovations and Δ is the suboptimality gap in Π . A significant hurdle of our formulation, which sets it aside from other online learning problems such as bandits, is that running a policy does not provide an unbiased estimate of its performance.

[When Expressivity Meets Trainability: Fewer than \$n\$ Neurons Can Work](#)

- Jiawei Zhang, Yushun Zhang, Mingyi Hong, Ruoyu Sun, Zhi-Quan Luo
- abstract@[open-review](#): Modern neural networks are often quite wide, causing large memory and computation costs. It is thus of great interest to train a narrower network. However, training narrow neural nets remains a challenging task. We ask two theoretical questions: Can narrow networks have as strong expressivity as wide ones? If so, does the loss function exhibit a benign optimization landscape? In this work, we provide partially affirmative answers to both questions for 1-hidden-layer networks with fewer than n (sample size) neurons when the activation is smooth. First, we prove that as long as the width $m \geq 2n/d$ (where d is the input dimension), its expressivity is strong, i.e., there exists at least one global minimizer with zero training loss. Second, we identify a nice local region with no local-min or saddle points. Nevertheless, it is not clear whether gradient descent can stay in this nice region. Third, we consider a constrained optimization formulation where the feasible region is the nice local region, and prove that every KKT point is a nearly global minimizer. It is expected that projected gradient methods converge to KKT points under mild technical conditions, but we leave the rigorous convergence analysis to future work. Thorough numerical results show that projected gradient methods on this constrained formulation significantly outperform SGD for training narrow neural nets.

[Analyzing the Confidentiality of Undistillable Teachers in Knowledge Distillation](#)

- Souvik Kundu, Qirui Sun, Yao Fu, Massoud Pedram, Peter Beerel
- abstract@[open-review](#): Knowledge distillation (KD) has recently been identified as a method that can unintentionally leak private information regarding the details of a teacher model to an unauthorized student. Recent research in developing undistillable nasty teachers that can protect model confidentiality has gained significant attention. However, the level of protection these nasty models offer has been largely untested. In this paper, we show that transferring knowledge to a shallow sub-section of a student can largely reduce a teacher's influence. By exploring the depth of the shallow subsection, we then present a distillation technique that enables a skeptical student model to learn even from a nasty teacher. To evaluate the efficacy of our skeptical students, we conducted experiments with several models with KD on both training data-available and data-free scenarios for various datasets. While distilling from nasty teachers, compared to the normal student models, skeptical students consistently provide superior classification performance of up to ~59.5%. Moreover, similar to normal students, skeptical students maintain high classification accuracy when distilled from a normal teacher, showing their efficacy irrespective of the teacher being nasty or not. We believe the ability of skeptical students to largely diminish the KD-immunity of potentially nasty teachers will motivate the research community to create more robust mechanisms for model confidentiality. We have open-sourced the code at <https://github.com/ksouvik52/Skeptical2021>

[High Probability Complexity Bounds for Line Search Based on Stochastic Oracles](#)

- Billy Jin, Katya Scheinberg, Miaolan Xie
- abstract@[open-review](#): We consider a line-search method for continuous optimization under a stochastic setting where the function values and gradients are available only through inexact probabilistic zeroth and first-order oracles. These oracles capture multiple standard settings including expected loss minimization and zeroth-order optimization. Moreover, our framework is very general and allows the function and gradient estimates to be biased. The proposed algorithm is simple to describe, easy to implement, and uses these oracles in a similar way as the standard deterministic line search uses exact function and gradient values. Under fairly general conditions on the oracles, we derive a high probability tail bound on the iteration complexity of the algorithm when applied to non-convex smooth functions. These results are stronger than those for other existing stochastic line search methods and apply in more general settings.

[Pay Attention to MLPs](#)

- Hanxiao Liu, Zihang Dai, David So, Quoc V Le
- abstract@[open-review](#): Transformers have become one of the most important architectural innovations in deep learning and have enabled many breakthroughs over the past few years. Here we propose a simple network architecture, gMLP, based solely on MLPs with gating, and show that it can perform as well as Transformers in key language and vision applications. Our comparisons show that self-attention is not critical for Vision Transformers, as gMLP can achieve the same accuracy. For BERT, our model achieves parity with Transformers on pretraining perplexity and is better on some downstream NLP tasks. On finetuning tasks where gMLP performs worse, making the gMLP model substantially larger can close the gap with Transformers. In general, our experiments show that gMLP can scale as well as Transformers over increased data and compute.

[An Image is Worth More Than a Thousand Words: Towards Disentanglement in The Wild](#)

- Aviv Gabbay, Niv Cohen, Yedid Hoshen
- abstract@[open-review](#): Unsupervised disentanglement has been shown to be theoretically impossible without inductive biases on the models and the data. As an alternative approach, recent methods rely on limited supervision to disentangle the factors of variation and allow their identifiability. While annotating the true generative factors is only required for a limited number of observations, we argue that it is infeasible to enumerate all the factors of variation that describe a real-world image distribution. To this end, we propose a method for disentangling a set of factors which are only partially labeled, as well as separating the complementary set of residual factors that are never explicitly specified. Our success in this challenging setting, demonstrated on synthetic benchmarks, gives rise to leveraging off-the-shelf image descriptors to partially annotate a subset of attributes in real image domains (e.g. of human faces) with minimal manual effort. Specifically, we use a recent language-image embedding model (CLIP) to annotate a set of attributes of interest in a zero-shot manner and demonstrate state-of-the-art disentangled image manipulation results.

[Dynamics of Stochastic Momentum Methods on Large-scale, Quadratic Models](#)

- Courtney Paquette, Elliot Paquette
- abstract@[open-review](#): We analyze a class of stochastic gradient algorithms with momentum on a high-dimensional random least squares problem. Our framework, inspired by random matrix theory, provides an exact (deterministic) characterization for the sequence of function values produced by these algorithms which is expressed only in terms of the eigenvalues of the Hessian. This leads to simple expressions for nearly-optimal hyperparameters, a description of the limiting neighborhood, and average-case complexity. As a consequence, we show that (small-batch) stochastic heavy-ball momentum with a fixed momentum parameter provides no actual performance improvement over SGD when step sizes are adjusted correctly. For contrast, in the non-strongly convex setting, it is possible to get a large improvement over SGD using momentum. By introducing hyperparameters that depend on the number of samples, we propose a new algorithm sDANA (stochastic dimension adjusted Nesterov acceleration) which obtains an asymptotically optimal average-case complexity while remaining linearly convergent in the strongly convex setting without adjusting parameters.

[Adversarial Examples in Multi-Layer Random ReLU Networks](#)

- Peter Bartlett, Sebastien Bubeck, Yeshwanth Cherapanamjeri

- abstract@[open-review](#): We consider the phenomenon of adversarial examples in ReLU networks with independent Gaussian parameters. For networks of constant depth and with a large range of widths (for instance, it suffices if the width of each layer is polynomial in that of any other layer), small perturbations of input vectors lead to large changes of outputs. This generalizes results of Daniely and Schacham (2020) for networks of rapidly decreasing width and of Bubeck et al (2021) for two-layer networks. Our proof shows that adversarial examples arise in these networks because the functions they compute are \emph{locally} very similar to random linear functions. Bottleneck layers play a key role: the minimal width up to some point in the network determines scales and sensitivities of mappings computed up to that point. The main result is for networks with constant depth, but we also show that some constraint on depth is necessary for a result of this kind, because there are suitably deep networks that, with constant probability, compute a function that is close to constant.

[Efficient Statistical Assessment of Neural Network Corruption Robustness](#)

- Karim TIT, Teddy Furon, Mathias ROUSSET
- abstract@[open-review](#): We quantify the robustness of a trained network to input uncertainties with a stochastic simulation inspired by the field of Statistical Reliability Engineering. The robustness assessment is cast as a statistical hypothesis test: the network is deemed as locally robust if the estimated probability of failure is lower than a critical level. The procedure is based on an Importance Splitting simulation generating samples of rare events. We derive theoretical guarantees that are non-asymptotic w.r.t. sample size. Experiments tackling large scale networks outline the efficiency of our method making a low number of calls to the network function.

[A Highly-Efficient Group Elastic Net Algorithm with an Application to Function-On-Scalar Regression](#)

- Tobia Boschi, Matthew Reimherr, Francesca Chiaromonte
- abstract@[open-review](#): Feature Selection and Functional Data Analysis are two dynamic areas of research, with important applications in the analysis of large and complex data sets. Straddling these two areas, we propose a new highly efficient algorithm to perform Group Elastic Net with application to function-on-scalar feature selection, where a functional response is modeled against a very large number of potential scalar predictors. First, we introduce a new algorithm to solve Group Elastic Net in ultra-high dimensional settings, which exploits the sparsity structure of the Augmented Lagrangian to greatly reduce the computational burden. Next, taking advantage of the properties of Functional Principal Components, we extend our algorithm to the function-on-scalar regression framework. We use simulations to demonstrate the CPU time gains afforded by our approach compared to its best existing competitors, and present an application to data from a Genome-Wide Association Study on childhood obesity.

[Hierarchical Clustering: \$\\$O\(1\)\\$\$ -Approximation for Well-Clustered Graphs](#)

- Bogdan-Adrian Manghiuc, He Sun
- abstract@[open-review](#): Hierarchical clustering studies a recursive partition of a data set into clusters of successively smaller size, and is a fundamental problem in data analysis. In this work we study the cost function for hierarchical clustering introduced by Dasgupta, and present two polynomial-time approximation algorithms: Our first result is an $\$O(1)\$$ -approximation algorithm for graphs of high conductance. Our simple construction bypasses complicated recursive routines of finding sparse cuts known in the literature. Our second and main result is an $\$O(1)\$$ -approximation algorithm for a wide family of graphs that exhibit a well-defined structure of clusters. This result generalises the previous state-of-the-art, which holds only for graphs generated from stochastic models. The significance of our work is demonstrated by the empirical analysis on both synthetic and real-world data sets, on which our presented algorithm outperforms the previously proposed algorithm for graphs with a well-defined cluster structure.

[Realistic evaluation of transductive few-shot learning](#)

- Olivier Veilleux, Malik Boudiaf, Pablo Piantanida, Ismail Ben Ayed
- abstract@[open-review](#): Transductive inference is widely used in few-shot learning, as it leverages the statistics of the unlabeled query set of a few-shot task, typically yielding substantially better performances than its inductive counterpart. The current few-shot benchmarks use perfectly class-balanced tasks at inference. We argue that such an artificial regularity is unrealistic, as it assumes that the marginal label probability of the testing samples is known and fixed to the uniform distribution. In fact, in realistic scenarios, the unlabeled query sets come with arbitrary and unknown label marginals. We introduce and study the effect of arbitrary class distributions within the query sets of few-shot tasks at inference, removing the class-balance artefact. Specifically, we model the marginal probabilities of the classes as Dirichlet-distributed random variables, which yields a principled and realistic sampling within the simplex. This leverages the current few-shot benchmarks, building testing tasks with arbitrary class distributions. We evaluate experimentally state-of-the-art transductive methods over 3 widely used data sets, and observe, surprisingly, substantial performance drops, even below inductive methods in some cases. Furthermore, we propose a generalization of the mutual-information loss, based on α -divergences, which can handle effectively class-distribution variations. Empirically, we show that our transductive α -divergence optimization outperforms state-of-the-art methods across several data sets, models and few-shot settings.

[Qu-ANTI-zation: Exploiting Quantization Artifacts for Achieving Adversarial Outcomes](#)

- Sanghyun Hong, Michael-Andrei Panaitescu-Liess, Yigitcan Kaya, Tudor Dumitras
- abstract@[open-review](#): Quantization is a popular technique that transforms the parameter representation of a neural network from floating-point numbers into lower-precision ones (e.g., 8-bit integers). It reduces the memory footprint and the computational cost at inference, facilitating the deployment of resource-hungry models. However, the parameter perturbations caused by this transformation result in behavioral disparities between the model before and after quantization. For example, a quantized model can misclassify some test-time samples that are otherwise classified correctly. It is not known whether such differences lead to a new security vulnerability. We hypothesize that an adversary may control this disparity to introduce specific behaviors that activate upon quantization. To study this hypothesis, we weaponize quantization-aware training and propose a new training framework to implement adversarial quantization outcomes. Following this framework, we present three attacks we carry out with quantization: (i) an indiscriminate attack for significant accuracy loss; (ii) a targeted attack against specific samples; and (iii) a backdoor attack for controlling the model with an input trigger. We further show that a single compromised model defeats multiple quantization schemes, including robust quantization techniques. Moreover, in a federated learning scenario, we demonstrate that a set of malicious participants who conspire can inject our quantization-activated backdoor. Lastly, we discuss potential counter-measures and show that only re-training consistently removes the attack artifacts. Our code is available at <https://github.com/Secure-AI-Systems-Group/Qu-ANTI-zation>

[Differentially Private Stochastic Optimization: New Results in Convex and Non-Convex Settings](#)

- Raef Bassily, Cristóbal Guzmán, Michael Menart
- abstract@[open-review](#): We study differentially private stochastic optimization in convex and non-convex settings. For the convex case, we focus on the family of non-smooth generalized linear losses (GLLs). Our algorithm for the $\$ell_2\$$ setting achieves optimal excess population risk in near-linear time, while the best known differentially private algorithms for general convex losses run in super-linear time. Our algorithm for the $\$ell_1\$$ setting has nearly-optimal excess population risk $\$tilde{O}\$ \big(\sqrt{\frac{\log d}{n}}\big)$, and circumvents the dimension dependent lower bound of $\$tilde{O}\$$ for general non-smooth convex losses. In the differentially private non-convex setting, we provide several new algorithms for approximating stationary points of the population risk. For the $\$ell_1\$$ -case with smooth losses and polyhedral constraint, we provide the first nearly dimension independent rate, $\$tilde{O}\$ \big(\frac{\log^{2/3} d}{n^{1/3}}\big)$ in linear time. For the constrained $\$ell_2\$$ -case, with smooth losses, we obtain a linear-time algorithm with rate $\$tilde{O}\$ \big(\frac{1}{d^{1/10}} + \frac{1}{d^{1/2}}\big)^{1/5}$. Finally, for the $\$ell_2\$$ -case we provide the first

method for $\{\text{em non-smooth weakly convex}\}$ stochastic optimization with rate $\tilde{O}(\frac{1}{n^{1/4}} + \frac{d}{n^2})^{1/6}$ which matches the best existing non-private algorithm when $d = O(\sqrt{n})$. We also extend all our results above for the non-convex ℓ_2 setting to the ℓ_p setting, where $1 < p \leq 2$, with only polylogarithmic (in the dimension) overhead in the rates.

[TacticZero: Learning to Prove Theorems from Scratch with Deep Reinforcement Learning](#)

- Minchao Wu, Michael Norrish, Christian Walder, Amir Dezfouli
- abstract@[open-review](#): We propose a novel approach to interactive theorem-proving (ITP) using deep reinforcement learning. The proposed framework is able to learn proof search strategies as well as tactic and arguments prediction in an end-to-end manner. We formulate the process of ITP as a Markov decision process (MDP) in which each state represents a set of potential derivation paths. This structure allows us to introduce a novel backtracking mechanism which enables the agent to efficiently discard (predicted) dead-end derivations and restart the derivation from promising alternatives. We implement the framework in the HOL theorem prover. Experimental results show that the framework using learned search strategies outperforms existing automated theorem provers (i.e., hammers) available in HOL when evaluated on unseen problems. We further elaborate the role of key components of the framework using ablation studies.

[Integrating Tree Path in Transformer for Code Representation](#)

- Han Peng, Ge Li, Wenhan Wang, YunFei Zhao, Zhi Jin
- abstract@[open-review](#): Learning distributed representation of source code requires modelling its syntax and semantics. Recent state-of-the-art models leverage highly structured source code representations, such as the syntax trees and paths therein. In this paper, we investigate two representative path encoding methods shown in previous research work and integrate them into the attention module of Transformer. We draw inspiration from the ideas of positional encoding and modify them to incorporate these path encoding. Specifically, we encode both the pairwise path between tokens of source code and the path from the leaf node to the tree root for each token in the syntax tree. We explore the interaction between these two kinds of paths by integrating them into the unified Transformer framework. The detailed empirical study for path encoding methods also leads to our novel state-of-the-art representation model TPTrans, which finally outperforms strong baselines. Extensive experiments and ablation studies on code summarization across four different languages demonstrate the effectiveness of our approaches. We release our code at <https://github.com/AwdHanPeng/TPTrans>.

[Twins: Revisiting the Design of Spatial Attention in Vision Transformers](#)

- Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, Chunhua Shen
- abstract@[open-review](#): Very recently, a variety of vision transformer architectures for dense prediction tasks have been proposed and they show that the design of spatial attention is critical to their success in these tasks. In this work, we revisit the design of the spatial attention and demonstrate that a carefully devised yet simple spatial attention mechanism performs favorably against the state-of-the-art schemes. As a result, we propose two vision transformer architectures, namely, Twins-PCPVT and Twins-SVT. Our proposed architectures are highly efficient and easy to implement, only involving matrix multiplications that are highly optimized in modern deep learning frameworks. More importantly, the proposed architectures achieve excellent performance on a wide range of visual tasks including image-level classification as well as dense detection and segmentation. The simplicity and strong performance suggest that our proposed architectures may serve as stronger backbones for many vision tasks.

[Evaluating State-of-the-Art Classification Models Against Bayes Optimality](#)

- Ryan Theisen, Huan Wang, Lav R. Varshney, Caiming Xiong, Richard Socher
- abstract@[open-review](#): Evaluating the inherent difficulty of a given data-driven classification problem is important for establishing absolute benchmarks and evaluating progress in the field. To this end, a natural quantity to consider is the Bayes error, which measures the optimal classification error theoretically achievable for a given data distribution. While generally an intractable quantity, we show that we can compute the exact Bayes error of generative models learned using normalizing flows. Our technique relies on a fundamental result, which states that the Bayes error is invariant under invertible transformation. Therefore, we can compute the exact Bayes error of the learned flow models by computing it for Gaussian base distributions, which can be done efficiently using Holmes-Diaconis-Ross integration. Moreover, we show that by varying the temperature of the learned flow models, we can generate synthetic datasets that closely resemble standard benchmark datasets, but with almost any desired Bayes error. We use our approach to conduct a thorough investigation of state-of-the-art classification models, and find that in some --- but not all --- cases, these models are capable of obtaining accuracy very near optimal. Finally, we use our method to evaluate the intrinsic "hardness" of standard benchmark datasets.

[Data-Efficient Instance Generation from Instance Discrimination](#)

- Ceyuan Yang, Yujun Shen, Yinghao Xu, Bolei Zhou
- abstract@[open-review](#): Generative Adversarial Networks (GANs) have significantly advanced image synthesis, however, the synthesis quality drops significantly given a limited amount of training data. To improve the data efficiency of GAN training, prior work typically employs data augmentation to mitigate the overfitting of the discriminator yet still learn the discriminator with a bi-classification (*i.e.*, real vs. fake) task. In this work, we propose a data-efficient Instance Generation (*InsGen*) method based on instance discrimination. Concretely, besides differentiating the real domain from the fake domain, the discriminator is required to distinguish every individual image, no matter it comes from the training set or from the generator. In this way, the discriminator can benefit from the infinite synthesized samples for training, alleviating the overfitting problem caused by insufficient training data. A noise perturbation strategy is further introduced to improve its discriminative power. Meanwhile, the learned instance discrimination capability from the discriminator is in turn exploited to encourage the generator for diverse generation. Extensive experiments demonstrate the effectiveness of our method on a variety of datasets and training settings. Noticeably, on the setting of 2K training images from the FFHQ dataset, we outperform the state-of-the-art approach with 23.5% FID improvement.

[Reliable Post hoc Explanations: Modeling Uncertainty in Explainability](#)

- Dylan Slack, Anna Hilgard, Sameer Singh, Himabindu Lakkaraju
- abstract@[open-review](#): As black box explanations are increasingly being employed to establish model credibility in high stakes settings, it is important to ensure that these explanations are accurate and reliable. However, prior work demonstrates that explanations generated by state-of-the-art techniques are inconsistent, unstable, and provide very little insight into their correctness and reliability. In addition, these methods are also computationally inefficient, and require significant hyper-parameter tuning. In this paper, we address the aforementioned challenges by developing a novel Bayesian framework for generating local explanations along with their associated uncertainty. We instantiate this framework to obtain Bayesian versions of LIME and KernelSHAP which output credible intervals for the feature importances, capturing the associated uncertainty. The resulting explanations not only enable us to make concrete inferences about their quality (*e.g.*, there is a 95% chance that the feature importance lies within the given range), but are also highly consistent and stable. We carry out a detailed theoretical analysis that leverages the aforementioned uncertainty to estimate how many perturbations to sample, and how to sample for faster convergence. This work makes the first attempt at addressing several critical issues with popular explanation methods in one shot, thereby generating consistent, stable, and reliable explanations with guarantees in a computationally efficient manner. Experimental evaluation with multiple real world datasets and user studies demonstrate that the efficacy of the proposed framework.

[Learning Graph Models for Retrosynthesis Prediction](#)

- Vignesh Ram Somnath, Charlotte Bunne, Connor Coley, Andreas Krause, Regina Barzilay
- abstract@[open-review](#): Retrosynthesis prediction is a fundamental problem in organic synthesis, where the task is to identify precursor molecules that can be used to synthesize a target molecule. A key consideration in building neural models for this task is aligning model design with strategies adopted by chemists. Building on this viewpoint, this paper introduces a graph-based approach that capitalizes on the idea that the graph topology of precursor molecules is largely unaltered during a chemical reaction. The model first predicts the set of graph edits transforming the target into incomplete molecules called synthons. Next, the model learns to expand synthons into complete molecules by attaching relevant leaving groups. This decomposition simplifies the architecture, making its predictions more interpretable, and also amenable to manual correction. Our model achieves a top-1 accuracy of 53.7%, outperforming previous template-free and semi-template-based methods.

[Differentiable Equilibrium Computation with Decision Diagrams for Stackelberg Models of Combinatorial Congestion Games](#)

- Shinsaku Sakaue, Kengo Nakamura
- abstract@[open-review](#): We address Stackelberg models of combinatorial congestion games (CCGs); we aim to optimize the parameters of CCGs so that the selfish behavior of non-atomic players attains desirable equilibria. This model is essential for designing such social infrastructures as traffic and communication networks. Nevertheless, computational approaches to the model have not been thoroughly studied due to two difficulties: (I) bilevel-programming structures and (II) the combinatorial nature of CCGs. We tackle them by carefully combining (I) the idea of \textit{differentiable} optimization and (II) data structures called \textit{zero-suppressed binary decision diagrams} (ZDDs), which can compactly represent sets of combinatorial strategies. Our algorithm numerically approximates the equilibria of CCGs, which we can differentiate with respect to parameters of CCGs by automatic differentiation. With the resulting derivatives, we can apply gradient-based methods to Stackelberg models of CCGs. Our method is tailored to induce Nesterov's acceleration and can fully utilize the empirical compactness of ZDDs. These technical advantages enable us to deal with CCGs with a vast number of combinatorial strategies. Experiments on real-world network design instances demonstrate the practicality of our method.

[Inverse Optimal Control Adapted to the Noise Characteristics of the Human Sensorimotor System](#)

- Matthias Schultheis, Dominik Straub, Constantin A. Rothkopf
- abstract@[open-review](#): Computational level explanations based on optimal feedback control with signal-dependent noise have been able to account for a vast array of phenomena in human sensorimotor behavior. However, commonly a cost function needs to be assumed for a task and the optimality of human behavior is evaluated by comparing observed and predicted trajectories. Here, we introduce inverse optimal control with signal-dependent noise, which allows inferring the cost function from observed behavior. To do so, we formalize the problem as a partially observable Markov decision process and distinguish between the agent's and the experimenter's inference problems. Specifically, we derive a probabilistic formulation of the evolution of states and belief states and an approximation to the propagation equation in the linear-quadratic Gaussian problem with signal-dependent noise. We extend the model to the case of partial observability of state variables from the point of view of the experimenter. We show the feasibility of the approach through validation on synthetic data and application to experimental data. Our approach enables recovering the costs and benefits implicit in human sequential sensorimotor behavior, thereby reconciling normative and descriptive approaches in a computational framework.

[Deep Neural Networks as Point Estimates for Deep Gaussian Processes](#)

- Vincent Dutordoir, James Hensman, Mark van der Wilk, Carl Henrik Ek, Zoubin Ghahramani, Nicolas Durrande
- abstract@[open-review](#): Neural networks and Gaussian processes are complementary in their strengths and weaknesses. Having a better understanding of their relationship comes with the promise to make each method benefit from the strengths of the other. In this work, we establish an equivalence between the forward passes of neural networks and (deep) sparse Gaussian process models. The theory we develop is based on interpreting activation functions as interdomain inducing features through a rigorous analysis of the interplay between activation functions and kernels. This results in models that can either be seen as neural networks with improved uncertainty prediction or deep Gaussian processes with increased prediction accuracy. These claims are supported by experimental results on regression and classification datasets.

[Locality defeats the curse of dimensionality in convolutional teacher-student scenarios](#)

- Alessandro Favero, Francesco Cagnetta, Matthieu Wyart
- abstract@[open-review](#): Convolutional neural networks perform a local and translationally-invariant treatment of the data: quantifying which of these two aspects is central to their success remains a challenge. We study this problem within a teacher-student framework for kernel regression, using 'convolutional' kernels inspired by the neural tangent kernel of simple convolutional architectures of given filter size. Using heuristic methods from physics, we find in the ridgeless case that locality is key in determining the learning curve exponent β (that relates the test error $\epsilon_{\text{test}} \sim P^{-\beta}$ to the size of the training set P), whereas translational invariance is not. In particular, if the filter size of the teacher s_t is smaller than that of the student s_s , β is a function of s_s only and does not depend on the input dimension. We confirm our predictions on β empirically. We conclude by proving, under a natural universality assumption, that performing kernel regression with a ridge that decreases with the size of the training set leads to similar learning curve exponents to those we obtain in the ridgeless case.

[Causal Identification with Matrix Equations](#)

- Sanghack Lee, Elias Bareinboim
- abstract@[open-review](#): Causal effect identification is concerned with determining whether a causal effect is computable from a combination of qualitative assumptions about the underlying system (e.g., a causal graph) and distributions collected from this system. Many identification algorithms exclusively rely on graphical criteria made of a non-trivial combination of probability axioms, do-calculus, and refined c-factorization (e.g., Lee & Bareinboim, 2020). In a sequence of increasingly sophisticated results, it has been shown how proxy variables can be used to identify certain effects that would not be otherwise recoverable in challenging scenarios through solving matrix equations (e.g., Kuroki & Pearl, 2014; Miao et al., 2018). In this paper, we develop a new causal identification algorithm which utilizes both graphical criteria and matrix equations. Specifically, we first characterize the relationships between certain graphically-driven formulae and matrix multiplications. With such characterizations, we broaden the spectrum of proxy variable based identification conditions and further propose novel intermediary criteria based on the pseudoinverse of a matrix. Finally, we devise a causal effect identification algorithm, which accepts as input a collection of marginal, conditional, and interventional distributions, integrating enriched matrix-based criteria into a graphical identification approach.

[Private and Non-private Uniformity Testing for Ranking Data](#)

- Róbert Busa-Fekete, Dimitris Fotakis, Emmanouil Zampetakis
- abstract@[open-review](#): We study the problem of uniformity testing for statistical data that consists of rankings over m items where the alternative class is restricted to Mallows models with single parameter. Testing ranking data is challenging because of the size of the large domain that is factorial in m , therefore the tester needs to take advantage of some structure of the alternative class. We show that uniform distribution can be distinguished from Mallows model with $O(m^{-1/2})$ samples based on simple pairwise statistics, which allows us to test uniformity using only two samples, if m is large enough. We also consider uniformity testing with central and locally differential private (DP) constraints. We present a central DP algorithm that requires $O(\max \{1/\epsilon_0, 1/\sqrt{m}\})$ where ϵ_0 is the privacy budget parameter. Interestingly, our uniformity testing algorithm is straightforward to apply in the local DP scenario by its nature, since it works with binary statistics that is extracted from the ranking data. We carry out large-scale experiments, including $m=10000$, to show that these testing algorithms scales very gracefully with the number of items.

Model-Based Reinforcement Learning via Imagination with Derived Memory

- Yao Mu, Yuzheng Zhuang, Bin Wang, Guangxiang Zhu, Wulong Liu, Jianyu Chen, Ping Luo, Shengbo Li, Chongjie Zhang, Jianye Hao
- abstract@[open-review](#): Model-based reinforcement learning aims to improve the sample efficiency of policy learning by modeling the dynamics of the environment. Recently, the latent dynamics model is further developed to enable fast planning in a compact space. It summarizes the high-dimensional experiences of an agent, which mimics the memory function of humans. Learning policies via imagination with the latent model shows great potential for solving complex tasks. However, only considering memories from the true experiences in the process of imagination could limit its advantages. Inspired by the memory prosthesis proposed by neuroscientists, we present a novel model-based reinforcement learning framework called Imagining with Derived Memory (IDM). It enables the agent to learn policy from enriched diverse imagination with prediction-reliability weight, thus improving sample efficiency and policy robustness. Experiments on various high-dimensional visual control tasks in the DMControl benchmark demonstrate that IDM outperforms previous state-of-the-art methods in terms of policy robustness and further improves the sample efficiency of the model-based method.

Compositional Transformers for Scene Generation

- Dor Arad Hudson, Larry Zitnick
- abstract@[open-review](#): We introduce the GANformer2 model, an iterative object-oriented transformer, explored for the task of generative modeling. The network incorporates strong and explicit structural priors, to reflect the compositional nature of visual scenes, and synthesizes images through a sequential process. It operates in two stages: a fast and lightweight planning phase, where we draft a high-level scene layout, followed by an attention-based execution phase, where the layout is being refined, evolving into a rich and detailed picture. Our model moves away from conventional black-box GAN architectures that feature a flat and monolithic latent space towards a transparent design that encourages efficiency, controllability and interpretability. We demonstrate GANformer2's strengths and qualities through a careful evaluation over a range of datasets, from multi-object CLEVR scenes to the challenging COCO images, showing it successfully achieves state-of-the-art performance in terms of visual quality, diversity and consistency. Further experiments demonstrate the model's disentanglement and provide a deeper insight into its generative process, as it proceeds step-by-step from a rough initial sketch, to a detailed layout that accounts for objects' depths and dependencies, and up to the final high-resolution depiction of vibrant and intricate real-world scenes. See <https://github.com/dorarad/gansformer> for model implementation.

An Exponential Lower Bound for Linearly Realizable MDP with Constant Suboptimality Gap

- Yuanhao Wang, Ruosong Wang, Sham Kakade
- abstract@[open-review](#): A fundamental question in the theory of reinforcement learning is: suppose the optimal Q -function lies in the linear span of a given d -dimensional feature mapping, is sample-efficient reinforcement learning (RL) possible? The recent and remarkable result of Weisz et al. (2020) resolves this question in the negative, providing an exponential (in d) sample size lower bound, which holds even if the agent has access to a generative model of the environment. One may hope that such a lower can be circumvented with an even stronger assumption that there is a constant gap between the optimal Q -value of the best action and that of the second-best action (for all states); indeed, the construction in Weisz et al. (2020) relies on having an exponentially small gap. This work resolves this subsequent question, showing that an exponential sample complexity lower bound still holds even if a constant gap is assumed. Perhaps surprisingly, this result implies an exponential separation between the online RL setting and the generative model setting, where sample-efficient RL is in fact possible in the latter setting with a constant gap. Complementing our negative hardness result, we give two positive results showing that provably sample-efficient RL is possible either under an additional low-variance assumption or under a novel hypercontractivity assumption.

Combating Noise: Semi-supervised Learning by Region Uncertainty Quantification

- Zhenyu Wang, Ya-Li Li, Ye Guo, Shengjin Wang
- abstract@[open-review](#): Semi-supervised learning aims to leverage a large amount of unlabeled data for performance boosting. Existing works primarily focus on image classification. In this paper, we delve into semi-supervised learning for object detection, where labeled data are more labor-intensive to collect. Current methods are easily distracted by noisy regions generated by pseudo labels. To combat the noisy labeling, we propose noise-resistant semi-supervised learning by quantifying the region uncertainty. We first investigate the adverse effects brought by different forms of noise associated with pseudo labels. Then we propose to quantify the uncertainty of regions by identifying the noise-resistant properties of regions over different strengths. By importing the region uncertainty quantification and promoting multi-peak probability distribution output, we introduce uncertainty into training and further achieve noise-resistant learning. Experiments on both PASCAL VOC and MS COCO demonstrate the extraordinary performance of our method.

Reducing the Covariate Shift by Mirror Samples in Cross Domain Alignment

- Yin Zhao, minquan wang, Longjun Cai
- abstract@[open-review](#): Eliminating the covariate shift cross domains is one of the common methods to deal with the issue of domain shift in visual unsupervised domain adaptation. However, current alignment methods, especially the prototype based or sample-level based methods neglect the structural properties of the underlying distribution and even break the condition of covariate shift. To relieve the limitations and conflicts, we introduce a novel concept named (virtual) mirror, which represents the equivalent sample in another domain. The equivalent sample pairs, named mirror pairs reflect the natural correspondence of the empirical distributions. Then a mirror loss, which aligns the mirror pairs cross domains, is constructed to enhance the alignment of the domains. The proposed method does not distort the internal structure of the underlying distribution. We also provide theoretical proof that the mirror samples and mirror loss have better asymptotic properties in reducing the domain shift. By applying the virtual mirror and mirror loss to the generic unsupervised domain adaptation model, we achieved consistently superior performance on several mainstream benchmarks.

Permutation-Invariant Variational Autoencoder for Graph-Level Representation Learning

- Robin Winter, Frank Noe, Djork-Arné Clevert
- abstract@[open-review](#): Recently, there has been great success in applying deep neural networks on graph structured data. Most work, however, focuses on either node- or graph-level supervised learning, such as node, link or graph classification or node-level unsupervised learning (e.g. node clustering). Despite its wide range of possible applications, graph-level unsupervised learning has not received much attention yet. This might be mainly attributed to the high representation complexity of graphs, which can be represented by $n!$ equivalent adjacency matrices, where n is the number of nodes. In this work we address this issue by proposing a permutation-invariant variational autoencoder for graph structured data. Our proposed model indirectly learns to match the node ordering of input and output graph, without imposing a particular node ordering or performing expensive graph matching. We demonstrate the effectiveness of our proposed model for graph reconstruction, generation and interpolation and evaluate the expressive power of extracted representations for downstream graph-level classification and regression.

Causal Abstractions of Neural Networks

- Atticus Geiger, Hanson Lu, Thomas Icard, Christopher Potts
- abstract@[open-review](#): Structural analysis methods (e.g., probing and feature attribution) are increasingly important tools for neural network analysis. We propose a new structural analysis method grounded in a formal theory of causal abstraction that provides rich characterizations of model-internal representations and their roles in input/output behavior. In this method, neural representations are aligned with variables in interpretable causal models,

and then interchange interventions are used to experimentally verify that the neural representations have the causal properties of their aligned variables. We apply this method in a case study to analyze neural models trained on Multiply Quantified Natural Language Inference (MQNLI) corpus, a highly complex NLI dataset that was constructed with a tree-structured natural logic causal model. We discover that a BERT-based model with state-of-the-art performance successfully realizes parts of the natural logic model's causal structure, whereas a simpler baseline model fails to show any such structure, demonstrating that neural representations encode the compositional structure of MQNLI examples.

[Conic Blackwell Algorithm: Parameter-Free Convex-Concave Saddle-Point Solving](#)

- Julien Grand-Clément, Christian Kroer
- abstract@[open-review](#): We develop new parameter-free and scale-free algorithms for solving convex-concave saddle-point problems. Our results are based on a new simple regret minimizer, the Conic Blackwell Algorithm $^+$ (CBA $^+$), which attains $\mathcal{O}(1/\sqrt{T})$ average regret. Intuitively, our approach generalizes to other decision sets of interest ideas from the Counterfactual Regret minimization (CFR $^+$) algorithm, which has very strong practical performance for solving sequential games on simplexes. We show how to implement CBA $^+$ for the simplex, ℓ_p norm balls, and ellipsoidal confidence regions in the simplex, and we present numerical experiments for solving matrix games and distributionally robust optimization problems. Our empirical results show that CBA $^+$ is a simple algorithm that outperforms state-of-the-art methods on synthetic data and real data instances, without the need for any choice of step sizes or other algorithmic parameters.

[3DP3: 3D Scene Perception via Probabilistic Programming](#)

- Nishad Gothiskar, Marco Cusumano-Towner, Ben Zinberg, Matin Ghavamizadeh, Falk Pollok, Austin Garrett, Josh Tenenbaum, Dan Gutfreund, Vikash Mansinghka
- abstract@[open-review](#): We present 3DP3, a framework for inverse graphics that uses inference in a structured generative model of objects, scenes, and images. 3DP3 uses (i) voxel models to represent the 3D shape of objects, (ii) hierarchical scene graphs to decompose scenes into objects and the contacts between them, and (iii) depth image likelihoods based on real-time graphics. Given an observed RGB-D image, 3DP3's inference algorithm infers the underlying latent 3D scene, including the object poses and a parsimonious joint parametrization of these poses, using fast bottom-up pose proposals, novel involutive MCMC updates of the scene graph structure, and, optionally, neural object detectors and pose estimators. We show that 3DP3 enables scene understanding that is aware of 3D shape, occlusion, and contact structure. Our results demonstrate that 3DP3 is more accurate at 6DoF object pose estimation from real images than deep learning baselines and shows better generalization to challenging scenes with novel viewpoints, contact, and partial observability.

[Novel Upper Bounds for the Constrained Most Probable Explanation Task](#)

- Tahrima Rahman, Sara Rouhani, Vibhav Gogate
- abstract@[open-review](#): We propose several schemes for upper bounding the optimal value of the constrained most probable explanation (CMPE) problem. Given a set of discrete random variables, two probabilistic graphical models defined over them and a real number q , this problem involves finding an assignment of values to all the variables such that the probability of the assignment is maximized according to the first model and is bounded by q w.r.t. the second model. In prior work, it was shown that CMPE is a unifying problem with several applications and special cases including the nearest assignment problem, the decision preserving most probable explanation task and robust estimation. It was also shown that CMPE is NP-hard even on tractable models such as bounded treewidth networks and is hard for integer linear programming methods because it includes a dense global constraint. The main idea in our approach is to simplify the problem via Lagrange relaxation and decomposition to yield either a knapsack problem or the unconstrained most probable explanation (MPE) problem, and then solving the two problems, respectively using specialized knapsack algorithms and mini-buckets based upper bounding schemes. We evaluate our proposed scheme along several dimensions including quality of the bounds and computation time required on various benchmark graphical models and how it can be used to find heuristic, near-optimal feasible solutions in an example application pertaining to robust estimation and adversarial attacks on classifiers.

[Why Spectral Normalization Stabilizes GANs: Analysis and Improvements](#)

- Zinan Lin, Vyas Sekar, Giulia Fanti
- abstract@[open-review](#): Spectral normalization (SN) is a widely-used technique for improving the stability and sample quality of Generative Adversarial Networks (GANs). However, current understanding of SN's efficacy is limited. In this work, we show that SN controls two important failure modes of GAN training: exploding and vanishing gradients. Our proofs illustrate a (perhaps unintentional) connection with the successful LeCun initialization. This connection helps to explain why the most popular implementation of SN for GANs requires no hyper-parameter tuning, whereas stricter implementations of SN have poor empirical performance out-of-the-box. Unlike LeCun initialization which only controls gradient vanishing at the beginning of training, SN preserves this property throughout training. Building on this theoretical understanding, we propose a new spectral normalization technique: Bidirectional Scaled Spectral Normalization (BSSN), which incorporates insights from later improvements to LeCun initialization: Xavier initialization and Kaiming initialization. Theoretically, we show that BSSN gives better gradient control than SN. Empirically, we demonstrate that it outperforms SN in sample quality and training stability on several benchmark datasets.

[\\$\text{\texttt{Implicit}}^2\\$: Implicit Layers for Implicit Representations](#)

- Zhichun Huang, Shaojie Bai, J. Zico Kolter
- abstract@[open-review](#): Recent research in deep learning has investigated two very different forms of "implicitness": implicit representations model high-frequency data such as images or 3D shapes directly via a low-dimensional neural network (often using e.g., sinusoidal bases or nonlinearities); implicit layers, in contrast, refer to techniques where the forward pass of a network is computed via non-linear dynamical systems, such as fixed-point or differential equation solutions, with the backward pass computed via the implicit function theorem. In this work, we demonstrate that these two seemingly orthogonal concepts are remarkably well-suited for each other. In particular, we show that by exploiting fixed-point implicit layer to model implicit representations, we can substantially improve upon the performance of the conventional explicit-layer-based approach. Additionally, as implicit representation networks are typically trained in large-batch settings, we propose to leverage the property of implicit layers to amortize the cost of fixed-point forward/backward passes over training steps -- thereby addressing one of the primary challenges with implicit layers (that many iterations are required for the black-box fixed-point solvers). We empirically evaluated our method on learning multiple implicit representations for images, videos and audios, showing that our $\text{\texttt{Implicit}}^2$ approach substantially improve upon existing models while being both faster to train and much more memory efficient.

[Mean-based Best Arm Identification in Stochastic Bandits under Reward Contamination](#)

- Arpan Mukherjee, Ali Tajer, Pin-Yu Chen, Payel Das
- abstract@[open-review](#): This paper investigates the problem of best arm identification in $\{s_l$ contaminated} stochastic multi-arm bandits. In this setting, the rewards obtained from any arm are replaced by samples from an adversarial model with probability λ . A fixed confidence (infinite-horizon) setting is considered, where the goal of the learner is to identify the arm with the largest mean. Owing to the adversarial contamination of the rewards, each arm's mean is only partially identifiable. This paper proposes two algorithms, a gap-based algorithm and one based on the successive elimination, for best arm identification in sub-Gaussian bandits. These algorithms involve mean estimates that achieve the optimal error guarantee on the deviation of the true mean from the estimate asymptotically. Furthermore, these algorithms asymptotically achieve the optimal sample complexity.

Specifically, for the gap-based algorithm, the sample complexity is asymptotically optimal up to constant factors, while for the successive elimination-based algorithm, it is optimal up to logarithmic factors. Finally, numerical experiments are provided to illustrate the gains of the algorithms compared to the existing baselines.

[MADE: Exploration via Maximizing Deviation from Explored Regions](#)

- Tianjun Zhang, Paria Rashidinejad, Jiantao Jiao, Yuandong Tian, Joseph E. Gonzalez, Stuart Russell
- abstract@[open-review](#): In online reinforcement learning (RL), efficient exploration remains particularly challenging in high-dimensional environments with sparse rewards. In low-dimensional environments, where tabular parameterization is possible, count-based upper confidence bound (UCB) exploration methods achieve minimax near-optimal rates. However, it remains unclear how to efficiently implement UCB in realistic RL tasks that involve non-linear function approximation. To address this, we propose a new exploration approach via maximizing the deviation of the occupancy of the next policy from the explored regions. We add this term as an adaptive regularizer to the standard RL objective to balance exploration vs. exploitation. We pair the new objective with a provably convergent algorithm, giving rise to a new intrinsic reward that adjusts existing bonuses. The proposed intrinsic reward is easy to implement and combine with other existing RL algorithms to conduct exploration. As a proof of concept, we evaluate the new intrinsic reward on tabular examples across a variety of model-based and model-free algorithms, showing improvements over count-only exploration strategies. When tested on navigation and locomotion tasks from MiniGrid and DeepMind Control Suite benchmarks, our approach significantly improves sample efficiency over state-of-the-art methods.

[Variational Automatic Curriculum Learning for Sparse-Reward Cooperative Multi-Agent Problems](#)

- Jiayu Chen, Yuanxin Zhang, Yuanfan Xu, Huimin Ma, Huazhong Yang, Jiaming Song, Yu Wang, Yi Wu
- abstract@[open-review](#): We introduce an automatic curriculum algorithm, Variational Automatic Curriculum Learning (VACL), for solving challenging goal-conditioned cooperative multi-agent reinforcement learning problems. We motivate our curriculum learning paradigm through a variational perspective, where the learning objective can be decomposed into two terms: task learning on the current curriculum, and curriculum update to a new task distribution. Local optimization over the second term suggests that the curriculum should gradually expand the training tasks from easy to hard. Our VACL algorithm implements this variational paradigm with two practical components, task expansion and entity curriculum, which produces a series of training tasks over both the task configurations as well as the number of entities in the task. Experiment results show that VACL solves a collection of sparse-reward problems with a large number of agents. Particularly, using a single desktop machine, VACL achieves 98% coverage rate with 100 agents in the simple-spread benchmark and reproduces the ramp-use behavior originally shown in OpenAI's hide-and-seek project.

[Align before Fuse: Vision and Language Representation Learning with Momentum Distillation](#)

- Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, Steven Chu Hong Hoi
- abstract@[open-review](#): Large-scale vision and language representation learning has shown promising improvements on various vision-language tasks. Most existing methods employ a transformer-based multimodal encoder to jointly model visual tokens (region-based image features) and word tokens. Because the visual tokens and word tokens are unaligned, it is challenging for the multimodal encoder to learn image-text interactions. In this paper, we introduce a contrastive loss to ALign the image and text representations BEfore Fusing (ALBEF) them through cross-modal attention, which enables more grounded vision and language representation learning. Unlike most existing methods, our method does not require bounding box annotations nor high-resolution images. In order to improve learning from noisy web data, we propose momentum distillation, a self-training method which learns from pseudo-targets produced by a momentum model. We provide a theoretical analysis of ALBEF from a mutual information maximization perspective, showing that different training tasks can be interpreted as different ways to generate views for an image-text pair. ALBEF achieves state-of-the-art performance on multiple downstream vision-language tasks. On image-text retrieval, ALBEF outperforms methods that are pre-trained on orders of magnitude larger datasets. On VQA and NLVR², ALBEF achieves absolute improvements of 2.37% and 3.84% compared to the state-of-the-art, while enjoying faster inference speed. Code and models are available at <https://github.com/salesforce/ALBEF>.

[Variational Model Inversion Attacks](#)

- Kuan-Chieh Wang, YAN FU, Ke Li, Ashish Khisti, Richard Zemel, Alireza Makhzani
- abstract@[open-review](#): Given the ubiquity of deep neural networks, it is important that these models do not reveal information about sensitive data that they have been trained on. In model inversion attacks, a malicious user attempts to recover the private dataset used to train a supervised neural network. A successful model inversion attack should generate realistic and diverse samples that accurately describe each of the classes in the private dataset. In this work, we provide a probabilistic interpretation of model inversion attacks, and formulate a variational objective that accounts for both diversity and accuracy. In order to optimize this variational objective, we choose a variational family defined in the code space of a deep generative model, trained on a public auxiliary dataset that shares some structural similarity with the target dataset. Empirically, our method substantially improves performance in terms of target attack accuracy, sample realism, and diversity on datasets of faces and chest X-ray images.

[Graph Neural Networks with Adaptive Residual](#)

- Xiaorui Liu, Jiayuan Ding, Wei Jin, Han Xu, Yao Ma, Zitao Liu, Jiliang Tang
- abstract@[open-review](#): Graph neural networks (GNNs) have shown the power in graph representation learning for numerous tasks. In this work, we discover an interesting phenomenon that although residual connections in the message passing of GNNs help improve the performance, they immensely amplify GNNs' vulnerability against abnormal node features. This is undesirable because in real-world applications, node features in graphs could often be abnormal such as being naturally noisy or adversarially manipulated. We analyze possible reasons to understand this phenomenon and aim to design GNNs with stronger resilience to abnormal features. Our understandings motivate us to propose and derive a simple, efficient, interpretable, and adaptive message passing scheme, leading to a novel GNN with Adaptive Residual, AirGNN. Extensive experiments under various abnormal feature scenarios demonstrate the effectiveness of the proposed algorithm.

[Efficient Active Learning for Gaussian Process Classification by Error Reduction](#)

- Guang Zhao, Edward Dougherty, Byung-Jun Yoon, Francis Alexander, Xiaoning Qian
- abstract@[open-review](#): Active learning sequentially selects the best instance for labeling by optimizing an acquisition function to enhance data/label efficiency. The selection can be either from a discrete instance set (pool-based scenario) or a continuous instance space (query synthesis scenario). In this work, we study both active learning scenarios for Gaussian Process Classification (GPC). The existing active learning strategies that maximize the Estimated Error Reduction (EER) aim at reducing the classification error after training with the new acquired instance in a one-step-look-ahead manner. The computation of EER-based acquisition functions is typically prohibitive as it requires retraining the GPC with every new query. Moreover, as the EER is not smooth, it can not be combined with gradient-based optimization techniques to efficiently explore the continuous instance space for query synthesis. To overcome these critical limitations, we develop computationally efficient algorithms for EER-based active learning with GPC. We derive the joint predictive distribution of label pairs as a one-dimensional integral, as a result of which the computation of the acquisition function avoids retraining the GPC for each query, remarkably reducing the computational overhead. We also derive the gradient chain rule to efficiently calculate the gradient of the acquisition function, which leads to the first query synthesis active learning algorithm implementing EER-based strategies. Our experiments clearly demonstrate the computational efficiency of the proposed algorithms. We also benchmark our algorithms on both synthetic and real-world datasets, which show superior performance in terms of sampling efficiency compared to the existing state-of-the-art algorithms.

[Non-Asymptotic Analysis for Two Time-scale TDC with General Smooth Function Approximation](#)

- Yue Wang, Shaofeng Zou, Yi Zhou
- abstract@[open-review](#): Temporal-difference learning with gradient correction (TDC) is a two time-scale algorithm for policy evaluation in reinforcement learning. This algorithm was initially proposed with linear function approximation, and was later extended to the one with general smooth function approximation. The asymptotic convergence for the on-policy setting with general smooth function approximation was established in [Bhatnagar et al., 2009], however, the non-asymptotic convergence analysis remains unsolved due to challenges in the non-linear and two-time-scale update structure, non-convex objective function and the projection onto a time-varying tangent plane. In this paper, we develop novel techniques to address the above challenges and explicitly characterize the non-asymptotic error bound for the general off-policy setting with i.i.d. or Markovian samples, and show that it converges as fast as $\mathcal{O}(1/\sqrt{T})$ (up to a factor of $\mathcal{O}(\log T)$). Our approach can be applied to a wide range of value-based reinforcement learning algorithms with general smooth function approximation.

[A Little Robustness Goes a Long Way: Leveraging Robust Features for Targeted Transfer Attacks](#)

- Jacob Springer, Melanie Mitchell, Garrett Kenyon
- abstract@[open-review](#): Adversarial examples for neural network image classifiers are known to be transferable: examples optimized to be misclassified by a source classifier are often misclassified as well by classifiers with different architectures. However, targeted adversarial examples—optimized to be classified as a chosen target class—tend to be less transferable between architectures. While prior research on constructing transferable targeted attacks has focused on improving the optimization procedure, in this work we examine the role of the source classifier. Here, we show that training the source classifier to be "slightly robust"—that is, robust to small-magnitude adversarial examples—substantially improves the transferability of class-targeted and representation-targeted adversarial attacks, even between architectures as different as convolutional neural networks and transformers. The results we present provide insight into the nature of adversarial examples as well as the mechanisms underlying so-called "robust" classifiers.

[TriBERT: Human-centric Audio-visual Representation Learning](#)

- Tanzila Rahman, Mengyu Yang, Leonid Sigal
- abstract@[open-review](#): The recent success of transformer models in language, such as BERT, has motivated the use of such architectures for multi-modal feature learning and tasks. However, most multi-modal variants (e.g., ViLBERT) have limited themselves to visual-linguistic data. Relatively few have explored its use in audio-visual modalities, and none, to our knowledge, illustrate them in the context of granular audio-visual detection or segmentation tasks such as sound source separation and localization. In this work, we introduce TriBERT -- a transformer-based architecture, inspired by ViLBERT, which enables contextual feature learning across three modalities: vision, pose, and audio, with the use of flexible co-attention. The use of pose keypoints is inspired by recent works that illustrate that such representations can significantly boost performance in many audio-visual scenarios where often one or more persons are responsible for the sound explicitly (e.g., talking) or implicitly (e.g., sound produced as a function of human manipulating an object). From a technical perspective, as part of the TriBERT architecture, we introduce a learned visual tokenization scheme based on spatial attention and leverage weak-supervision to allow granular cross-modal interactions for visual and pose modalities. Further, we supplement learning with sound-source separation loss formulated across all three streams. We pre-train our model on the large MUSIC21 dataset and demonstrate improved performance in audio-visual sound source separation on that dataset as well as other datasets through fine-tuning. In addition, we show that the learned TriBERT representations are generic and significantly improve performance on other audio-visual tasks such as cross-modal audio-visual-pose retrieval by as much as 66.7% in top-1 accuracy.

[How does a Neural Network's Architecture Impact its Robustness to Noisy Labels?](#)

- Jingling Li, Mozhi Zhang, Keyulu Xu, John Dickerson, Jimmy Ba
- abstract@[open-review](#): Noisy labels are inevitable in large real-world datasets. In this work, we explore an area understudied by previous works --- how the network's architecture impacts its robustness to noisy labels. We provide a formal framework connecting the robustness of a network to the alignments between its architecture and target/noise functions. Our framework measures a network's robustness via the predictive power in its representations --- the test performance of a linear model trained on the learned representations using a small set of clean labels. We hypothesize that a network is more robust to noisy labels if its architecture is more aligned with the target function than the noise. To support our hypothesis, we provide both theoretical and empirical evidence across various neural network architectures and different domains. We also find that when the network is well-aligned with the target function, its predictive power in representations could improve upon state-of-the-art (SOTA) noisy-label-training methods in terms of test accuracy and even outperform sophisticated methods that use clean labels.

[Calibration and Consistency of Adversarial Surrogate Losses](#)

- Pranjal Awasthi, Natalie Frank, Anqi Mao, Mehryar Mohri, Yutao Zhong
- abstract@[open-review](#): Adversarial robustness is an increasingly critical property of classifiers in applications. The design of robust algorithms relies on surrogate losses since the optimization of the adversarial loss with most hypothesis sets is NP-hard. But, which surrogate losses should be used and when do they benefit from theoretical guarantees? We present an extensive study of this question, including a detailed analysis of the \mathcal{H} -calibration and \mathcal{H} -consistency of adversarial surrogate losses. We show that convex loss functions, or the supremum-based convex losses often used in applications, are not \mathcal{H} -calibrated for common hypothesis sets used in machine learning. We then give a characterization of \mathcal{H} -calibration and prove that some surrogate losses are indeed \mathcal{H} -calibrated for the adversarial zero-one loss, with common hypothesis sets. In particular, we fix some calibration results presented in prior work for a family of linear models and significantly generalize the results to the nonlinear hypothesis sets. Next, we show that \mathcal{H} -calibration is not sufficient to guarantee consistency and prove that, in the absence of any distributional assumption, no continuous surrogate loss is consistent in the adversarial setting. This, in particular, proves that a claim made in prior work is inaccurate. Next, we identify natural conditions under which some surrogate losses that we describe in detail are \mathcal{H} -consistent. We also report a series of empirical results which show that many \mathcal{H} -calibrated surrogate losses are indeed not \mathcal{H} -consistent, and validate our theoretical assumptions. Our adversarial \mathcal{H} -consistency results are novel, even for the case where \mathcal{H} is the family of all measurable functions.

[The Value of Information When Deciding What to Learn](#)

- Dilip Arumugam, Benjamin Van Roy
- abstract@[open-review](#): All sequential decision-making agents explore so as to acquire knowledge about a particular target. It is often the responsibility of the agent designer to construct this target which, in rich and complex environments, constitutes an onerous burden; without full knowledge of the environment itself, a designer may forge a sub-optimal learning target that poorly balances the amount of information an agent must acquire to identify the target against the target's associated performance shortfall. While recent work has developed a connection between learning targets and rate-distortion theory to address this challenge and empower agents that decide what to learn in an automated fashion, the proposed algorithm does not optimally tackle the equally important challenge of efficient information acquisition. In this work, building upon the seminal design principle of information-directed sampling (Russo & Van Roy, 2014), we address this shortcoming directly to couple optimal information acquisition with the optimal design of learning targets. Along the way, we offer new insights into learning targets from the literature on rate-distortion theory before turning to empirical results that confirm the value of information when deciding what to learn.

Co-Adaptation of Algorithmic and Implementational Innovations in Inference-based Deep Reinforcement Learning

- Hiroki Furuta, Tadashi Kozuno, Tatsuya Matsushima, Yutaka Matsuo, Shixiang (Shane) Gu
- abstract@[open-review](#): Recently many algorithms were devised for reinforcement learning (RL) with function approximation. While they have clear algorithmic distinctions, they also have many implementation differences that are algorithm-independent and sometimes under-emphasized. Such mixing of algorithmic novelty and implementation craftsmanship makes rigorous analyses of the sources of performance improvements across algorithms difficult. In this work, we focus on a series of off-policy inference-based actor-critic algorithms -- MPO, AWR, and SAC -- to decouple their algorithmic innovations and implementation decisions. We present unified derivations through a single control-as-inference objective, where we can categorize each algorithm as based on either Expectation-Maximization (EM) or direct Kullback-Leibler (KL) divergence minimization and treat the rest of specifications as implementation details. We performed extensive ablation studies, and identified substantial performance drops whenever implementation details are mismatched for algorithmic choices. These results show which implementation or code details are co-adapted and co-evolved with algorithms, and which are transferable across algorithms: as examples, we identified that tanh Gaussian policy and network sizes are highly adapted to algorithmic types, while layer normalization and ELU are critical for MPO's performances but also transfer to noticeable gains in SAC. We hope our work can inspire future work to further demystify sources of performance improvements across multiple algorithms and allow researchers to build on one another's both algorithmic and implementational innovations.

Can fMRI reveal the representation of syntactic structure in the brain?

- Aniketh Janardhan Reddy, Leila Wehbe
- abstract@[open-review](#): While studying semantics in the brain, neuroscientists use two approaches. One is to identify areas that are correlated with semantic processing load. Another is to find areas that are predicted by the semantic representation of the stimulus words. However, most studies of syntax have focused only on identifying areas correlated with syntactic processing load. One possible reason for this discrepancy is that representing syntactic structure in an embedding space such that it can be used to model brain activity is a non-trivial computational problem. Another possible reason is that it is unclear if the low signal-to-noise ratio of neuroimaging tools such as functional Magnetic Resonance Imaging (fMRI) can allow us to reveal the correlates of complex (and perhaps subtle) syntactic representations. In this study, we propose novel multi-dimensional features that encode information about the syntactic structure of sentences. Using these features and fMRI recordings of participants reading a natural text, we model the brain representation of syntax. First, we find that our syntactic structure-based features explain additional variance in the brain activity of various parts of the language system, even after controlling for complexity metrics that capture processing load. At the same time, we see that regions well-predicted by syntactic features are distributed in the language system and are not distinguishable from those processing semantics. Our code and data will be available at <https://github.com/anikethjr/brainsyntacticrepresentations>.

Robust Implicit Networks via Non-Euclidean Contractions

- Saber Jafarpour, Alexander Davydov, Anton Proskurnikov, Francesco Bullo
- abstract@[open-review](#): Implicit neural networks, a.k.a., deep equilibrium networks, are a class of implicit-depth learning models where function evaluation is performed by solving a fixed point equation. They generalize classic feedforward models and are equivalent to infinite-depth weight-tied feedforward networks. While implicit models show improved accuracy and significant reduction in memory consumption, they can suffer from ill-posedness and convergence instability. This paper provides a new framework, which we call Non-Euclidean Monotone Operator Network (NEMON), to design well-posed and robust implicit neural networks based upon contraction theory for the non-Euclidean norm $\|\cdot\|_{\text{infty}}$. Our framework includes (i) a novel condition for well-posedness based on one-sided Lipschitz constants, (ii) an average iteration for computing fixed-points, and (iii) explicit estimates on input-output Lipschitz constants. Additionally, we design a training problem with the well-posedness condition and the average iteration as constraints and, to achieve robust models, with the input-output Lipschitz constant as a regularizer. Our $\|\cdot\|_{\text{infty}}$ well-posedness condition leads to a larger polytopic training search space than existing conditions and our average iteration enjoys accelerated convergence. Finally, we evaluate our framework in image classification through the MNIST and the CIFAR-10 datasets. Our numerical results demonstrate improved accuracy and robustness of the implicit models with smaller input-output Lipschitz bounds. Code is available at https://github.com/davydovalexander/Non-Euclidean_Mon_Op_Net.

A Kernel-based Test of Independence for Cluster-correlated Data

- Hongjiao Liu, Anna Plantinga, Yunhua Xiang, Michael Wu
- abstract@[open-review](#): The Hilbert-Schmidt Independence Criterion (HSIC) is a powerful kernel-based statistic for assessing the generalized dependence between two multivariate variables. However, independence testing based on the HSIC is not directly possible for cluster-correlated data. Such a correlation pattern among the observations arises in many practical situations, e.g., family-based and longitudinal data, and requires proper accommodation. Therefore, we propose a novel HSIC-based independence test to evaluate the dependence between two multivariate variables based on cluster-correlated data. Using the previously proposed empirical HSIC as our test statistic, we derive its asymptotic distribution under the null hypothesis of independence between the two variables but in the presence of sample correlation. Based on both simulation studies and real data analysis, we show that, with clustered data, our approach effectively controls type I error and has a higher statistical power than competing methods.

Efficient methods for Gaussian Markov random fields under sparse linear constraints

- David Bolin, Jonas Wallin
- abstract@[open-review](#): Methods for inference and simulation of linearly constrained Gaussian Markov Random Fields (GMRF) are computationally prohibitive when the number of constraints is large. In some cases, such as for intrinsic GMRFs, they may even be unfeasible. We propose a new class of methods to overcome these challenges in the common case of sparse constraints, where one has a large number of constraints and each only involves a few elements. Our methods rely on a basis transformation into blocks of constrained versus non-constrained subspaces, and we show that the methods greatly outperform existing alternatives in terms of computational cost. By combining the proposed methods with the stochastic partial differential equation approach for Gaussian random fields, we also show how to formulate Gaussian process regression with linear constraints in a GMRF setting to reduce computational cost. This is illustrated in two applications with simulated data.

Sparse is Enough in Scaling Transformers

- Sebastian Jaszczyk, Aakanksha Chowdhery, Afroz Mohiuddin, LUKASZ KAISER, Wojciech Gajewski, Henryk Michalewski, Jonni Kanerva
- abstract@[open-review](#): Large Transformer models yield impressive results on many tasks, but are expensive to train, or even fine-tune, and so slow at decoding that their use and study becomes out of reach. We address this problem by leveraging sparsity. We study sparse variants for all layers in the Transformer and propose Scaling Transformers, a family of next generation Transformer models that use sparse layers to scale efficiently and perform unbatched decoding much faster than the standard Transformer as we scale up the model size. Surprisingly, the sparse layers are enough to obtain the same perplexity as the standard Transformer with the same number of parameters. We also integrate with prior sparsity approaches to attention and enable fast inference on long sequences even with limited memory. This results in performance competitive to the state-of-the-art on long text summarization.

Sparse Training via Boosting Pruning Plasticity with Neuroregeneration

- Shiwei Liu, Tianlong Chen, Xiaohan Chen, Zahra Atashgahi, Lu Yin, Huanyu Kou, Li Shen, Mykola Pechenizkiy, Zhangyang Wang, Decebal Constantin Mocanu
- abstract@[open-review](#): Works on lottery ticket hypothesis (LTH) and single-shot network pruning (SNIP) have raised a lot of attention currently on post-training pruning (iterative magnitude pruning), and before-training pruning (pruning at initialization). The former method suffers from an extremely large computation cost and the latter usually struggles with insufficient performance. In comparison, during-training pruning, a class of pruning methods that simultaneously enjoys the training/inference efficiency and the comparable performance, temporarily, has been less explored. To better understand during-training pruning, we quantitatively study the effect of pruning throughout training from the perspective of pruning plasticity (the ability of the pruned networks to recover the original performance). Pruning plasticity can help explain several other empirical observations about neural network pruning in literature. We further find that pruning plasticity can be substantially improved by injecting a brain-inspired mechanism called neuroregeneration, i.e., to regenerate the same number of connections as pruned. We design a novel gradual magnitude pruning (GMP) method, named gradual pruning with zero-cost neuroregeneration (GraNet), that advances state of the art. Perhaps most impressively, its sparse-to-sparse version for the first time boosts the sparse-to-sparse training performance over various dense-to-sparse methods with ResNet-50 on ImageNet without extending the training time. We release all codes in <https://github.com/Shiweiliuiiiiiii/GraNet>.

[Low-Fidelity Video Encoder Optimization for Temporal Action Localization](#)

- Mengmeng Xu, Juan Manuel Perez Rua, Xiatian Zhu, Bernard Ghanem, Brais Martinez
- abstract@[open-review](#): Most existing temporal action localization (TAL) methods rely on a transfer learning pipeline: by first optimizing a video encoder on a large action classification dataset (i.e., source domain), followed by freezing the encoder and training a TAL head on the action localization dataset (i.e., target domain). This results in a task discrepancy problem for the video encoder – trained for action classification, but used for TAL. Intuitively, joint optimization with both the video encoder and TAL head is a strong baseline solution to this discrepancy. However, this is not operable for TAL subject to the GPU memory constraints, due to the prohibitive computational cost in processing long untrimmed videos. In this paper, we resolve this challenge by introducing a novel low-fidelity (LoFi) video encoder optimization method. Instead of always using the full training configurations in TAL learning, we propose to reduce the mini-batch composition in terms of temporal, spatial, or spatio-temporal resolution so that jointly optimizing the video encoder and TAL head becomes operable under the same memory conditions of a mid-range hardware budget. Crucially, this enables the gradients to flow backwards through the video encoder conditioned on a TAL supervision loss, favourably solving the task discrepancy problem and providing more effective feature representations. Extensive experiments show that the proposed LoFi optimization approach can significantly enhance the performance of existing TAL methods. Encouragingly, even with a lightweight ResNet18 based video encoder in a single RGB stream, our method surpasses two-stream (RGB + optical-flow) ResNet50 based alternatives, often by a good margin. Our code is publicly available at <https://github.com/saic-fi/lofactionlocalization>.

[On Provable Benefits of Depth in Training Graph Convolutional Networks](#)

- Weilin Cong, Morteza Ramezani, Mehrdad Mahdavi
- abstract@[open-review](#): Graph Convolutional Networks (GCNs) are known to suffer from performance degradation as the number of layers increases, which is usually attributed to over-smoothing. Despite the apparent consensus, we observe that there exists a discrepancy between the theoretical understanding of over-smoothing and the practical capabilities of GCNs. Specifically, we argue that over-smoothing does not necessarily happen in practice, a deeper model is provably expressive, can converge to global optimum with linear convergence rate, and achieve very high training accuracy as long as properly trained. Despite being capable of achieving high training accuracy, empirical results show that the deeper models generalize poorly on the testing stage and existing theoretical understanding of such behavior remains elusive. To achieve better understanding, we carefully analyze the generalization capability of GCNs, and show that the training strategies to achieve high training accuracy significantly deteriorate the generalization capability of GCNs. Motivated by these findings, we propose a decoupled structure for GCNs that detaches weight matrices from feature propagation to preserve the expressive power and ensure good generalization performance. We conduct empirical evaluations on various synthetic and real-world datasets to validate the correctness of our theory.

[Practical Near Neighbor Search via Group Testing](#)

- Joshua Engels, Benjamin Coleman, Anshumali Shrivastava
- abstract@[open-review](#): We present a new algorithm for the approximate near neighbor problem that combines classical ideas from group testing with locality-sensitive hashing (LSH). We reduce the near neighbor search problem to a group testing problem by designating neighbors as "positives," non-neighbors as "negatives," and approximate membership queries as group tests. We instantiate this framework using distance-sensitive Bloom Filters to Identify Near-Neighbor Groups (FLINNG). We prove that FLINNG has sub-linear query time and show that our algorithm comes with a variety of practical advantages. For example, FLINNG can be constructed in a single pass through the data, consists entirely of efficient integer operations, and does not require any distance computations. We conduct large-scale experiments on high-dimensional search tasks such as genome search, URL similarity search, and embedding search over the massive YFCC100M dataset. In our comparison with leading algorithms such as HNSW and FAISS, we find that FLINNG can provide up to a 10x query speedup with substantially smaller indexing time and memory.

[Baby Intuitions Benchmark \(BIB\): Discerning the goals, preferences, and actions of others](#)

- Kanishk Gandhi, Gala Stojnic, Brenden M. Lake, Moira R Dillon
- abstract@[open-review](#): To achieve human-like common sense about everyday life, machine learning systems must understand and reason about the goals, preferences, and actions of other agents in the environment. By the end of their first year of life, human infants intuitively achieve such common sense, and these cognitive achievements lay the foundation for humans' rich and complex understanding of the mental states of others. Can machines achieve generalizable, commonsense reasoning about other agents like human infants? The Baby Intuitions Benchmark (BIB) challenges machines to predict the plausibility of an agent's behavior based on the underlying causes of its actions. Because BIB's content and paradigm are adopted from developmental cognitive science, BIB allows for direct comparison between human and machine performance. Nevertheless, recently proposed, deep-learning-based agency reasoning models fail to show infant-like reasoning, leaving BIB an open challenge.

[Neural Hybrid Automata: Learning Dynamics With Multiple Modes and Stochastic Transitions](#)

- Michael Poli, Stefano Massaroli, Luca Scimeca, Sanghyuk Chun, Seong Joon Oh, Atsushi Yamashita, Hajime Asama, Jinkyoo Park, Animesh Garg
- abstract@[open-review](#): Effective control and prediction of dynamical systems require appropriate handling of continuous-time and discrete, event-triggered processes. Stochastic hybrid systems (SHSs), common across engineering domains, provide a formalism for dynamical systems subject to discrete, possibly stochastic, state jumps and multi-modal continuous-time flows. Despite the versatility and importance of SHSs across applications, a general procedure for the explicit learning of both discrete events and multi-mode continuous dynamics remains an open problem. This work introduces Neural Hybrid Automata (NHAs), a recipe for learning SHS dynamics without a priori knowledge on the number, mode parameters, and inter-modal transition dynamics. NHAs provide a systematic inference method based on normalizing flows, neural differential equations, and self-supervision. We showcase NHAs on several tasks, including mode recovery and flow learning in systems with stochastic transitions, and end-to-end learning of hierarchical robot controllers.

[Fast Projection onto the Capped Simplex with Applications to Sparse Regression in Bioinformatics](#)

- Man Shun Ang, Jianzhu Ma, Nianjun Liu, Kun Huang, Yijie Wang

- abstract@[open-review](#): We consider the problem of projecting a vector onto the so-called k-capped simplex, which is a hyper-cube cut by a hyperplane. For an n-dimensional input vector with bounded elements, we found that a simple algorithm based on Newton's method is able to solve the projection problem to high precision with a complexity roughly about $O(n)$, which has a much lower computational cost compared with the existing sorting-based methods proposed in the literature. We provide a theory for partial explanation and justification of the method. We demonstrate that the proposed algorithm can produce a solution of the projection problem with high precision on large scale datasets, and the algorithm is able to significantly outperform the state-of-the-art methods in terms of runtime (about 6-8 times faster than a commercial software with respect to CPU time for input vector with 1 million variables or more). We further illustrate the effectiveness of the proposed algorithm on solving sparse regression in a bioinformatics problem. Empirical results on the GWAS dataset (with 1,500,000 single-nucleotide polymorphisms) show that, when using the proposed method to accelerate the Projected Quasi-Newton (PQN) method, the accelerated PQN algorithm is able to handle huge-scale regression problem and it is more efficient (about 3-6 times faster) than the current state-of-the-art methods.

[The Many Faces of Adversarial Risk](#)

- Muni Sreenivas Pydi, Varun Jog
- abstract@[open-review](#): Adversarial risk quantifies the performance of classifiers on adversarially perturbed data. Numerous definitions of adversarial risk--not all mathematically rigorous and differing subtly in the details--have appeared in the literature. In this paper, we revisit these definitions, make them rigorous, and critically examine their similarities and differences. Our technical tools derive from optimal transport, robust statistics, functional analysis, and game theory. Our contributions include the following: generalizing Strassen's theorem to the unbalanced optimal transport setting with applications to adversarial classification with unequal priors; showing an equivalence between adversarial robustness and robust hypothesis testing with $\$infty$ -Wasserstein uncertainty sets; proving the existence of a pure Nash equilibrium in the two-player game between the adversary and the algorithm; and characterizing adversarial risk by the minimum Bayes error between distributions belonging to the $\$infty$ -Wasserstein uncertainty sets. Our results generalize and deepen recently discovered connections between optimal transport and adversarial robustness and reveal new connections to Choquet capacities and game theory.

[Meta-Adaptive Nonlinear Control: Theory and Algorithms](#)

- Guanya Shi, Kamyar Azizzadenesheli, Michael O'Connell, Soon-Jo Chung, Yisong Yue
- abstract@[open-review](#): We present an online multi-task learning approach for adaptive nonlinear control, which we call Online Meta-Adaptive Control (OMAC). The goal is to control a nonlinear system subject to adversarial disturbance and unknown \emph{environment-dependent} nonlinear dynamics, under the assumption that the environment-dependent dynamics can be well captured with some shared representation. Our approach is motivated by robot control, where a robotic system encounters a sequence of new environmental conditions that it must quickly adapt to. A key emphasis is to integrate online representation learning with established methods from control theory, in order to arrive at a unified framework that yields both control-theoretic and learning-theoretic guarantees. We provide instantiations of our approach under varying conditions, leading to the first non-asymptotic end-to-end convergence guarantee for multi-task nonlinear control. OMAC can also be integrated with deep representation learning. Experiments show that OMAC significantly outperforms conventional adaptive control approaches which do not learn the shared representation, in inverted pendulum and 6-DoF drone control tasks under varying wind conditions.

[Compositional Reinforcement Learning from Logical Specifications](#)

- Kishor Jothimurugan, Suguman Bansal, Osbert Bastani, Rajeev Alur
- abstract@[open-review](#): We study the problem of learning control policies for complex tasks given by logical specifications. Recent approaches automatically generate a reward function from a given specification and use a suitable reinforcement learning algorithm to learn a policy that maximizes the expected reward. These approaches, however, scale poorly to complex tasks that require high-level planning. In this work, we develop a compositional learning approach, called DIRL, that interleaves high-level planning and reinforcement learning. First, DIRL encodes the specification as an abstract graph; intuitively, vertices and edges of the graph correspond to regions of the state space and simpler sub-tasks, respectively. Our approach then incorporates reinforcement learning to learn neural network policies for each edge (sub-task) within a Dijkstra-style planning algorithm to compute a high-level plan in the graph. An evaluation of the proposed approach on a set of challenging control benchmarks with continuous state and action spaces demonstrates that it outperforms state-of-the-art baselines.

[Differentiable Quality Diversity](#)

- Matthew Fontaine, Stefanos Nikolaidis
- abstract@[open-review](#): Quality diversity (QD) is a growing branch of stochastic optimization research that studies the problem of generating an archive of solutions that maximize a given objective function but are also diverse with respect to a set of specified measure functions. However, even when these functions are differentiable, QD algorithms treat them as "black boxes", ignoring gradient information. We present the differentiable quality diversity (DQD) problem, a special case of QD, where both the objective and measure functions are first order differentiable. We then present MAP-Elites via a Gradient Arborescence (MEGA), a DQD algorithm that leverages gradient information to efficiently explore the joint range of the objective and measure functions. Results in two QD benchmark domains and in searching the latent space of a StyleGAN show that MEGA significantly outperforms state-of-the-art QD algorithms, highlighting DQD's promise for efficient quality diversity optimization when gradient information is available. Source code is available at <https://github.com/icaros-usc/dqd>.

[Credit Assignment Through Broadcasting a Global Error Vector](#)

- David Clark, L F Abbott, Sueyeon Chung
- abstract@[open-review](#): Backpropagation (BP) uses detailed, unit-specific feedback to train deep neural networks (DNNs) with remarkable success. That biological neural circuits appear to perform credit assignment, but cannot implement BP, implies the existence of other powerful learning algorithms. Here, we explore the extent to which a globally broadcast learning signal, coupled with local weight updates, enables training of DNNs. We present both a learning rule, called global error-vector broadcasting (GEVB), and a class of DNNs, called vectorized nonnegative networks (VNNs), in which this learning rule operates. VNNs have vector-valued units and nonnegative weights past the first layer. The GEVB learning rule generalizes three-factor Hebbian learning, updating each weight by an amount proportional to the inner product of the presynaptic activation and a globally broadcast error vector when the postsynaptic unit is active. We prove that these weight updates are matched in sign to the gradient, enabling accurate credit assignment. Moreover, at initialization, these updates are exactly proportional to the gradient in the limit of infinite network width. GEVB matches the performance of BP in VNNs, and in some cases outperforms direct feedback alignment (DFA) applied in conventional networks. Unlike DFA, GEVB successfully trains convolutional layers. Altogether, our theoretical and empirical results point to a surprisingly powerful role for a global learning signal in training DNNs.

[An Online Method for A Class of Distributionally Robust Optimization with Non-convex Objectives](#)

- Qi Qi, Zhishuai Guo, Yi Xu, Rong Jin, Tianbao Yang
- abstract@[open-review](#): In this paper, we propose a practical online method for solving a class of distributional robust optimization (DRO) with non-convex objectives, which has important applications in machine learning for improving the robustness of neural networks. In the literature, most methods for solving DRO are based on stochastic primal-dual methods. However, primal-dual methods for DRO suffer from several drawbacks: (1) manipulating a high-dimensional dual variable corresponding to the size of data is time expensive; (2) they are not friendly to online learning where data is coming

sequentially. To address these issues, we consider a class of DRO with an KL divergence regularization on the dual variables, transform the min-max problem into a compositional minimization problem, and propose practical duality-free online stochastic methods without requiring a large mini-batch size. We establish the state-of-the-art complexities of the proposed methods with and without a Polyak-Łojasiewicz (PL) condition of the objective. Empirical studies on large-scale deep learning tasks (i) demonstrate that our method can speed up the training by more than 2 times than baseline methods and save days of training time on a large-scale dataset with $\sim 265K$ images, and (ii) verify the supreme performance of DRO over Empirical Risk Minimization (ERM) on imbalanced datasets. Of independent interest, the proposed method can be also used for solving a family of stochastic compositional problems with state-of-the-art complexities.

[A single gradient step finds adversarial examples on random two-layers neural networks](#)

- Sébastien Bubeck, Yeshwanth Cherapanamjeri, Gauthier Gidel, Remi Tachet des Combes
- abstract@[open-review](#): Daniely and Schacham recently showed that gradient descent finds adversarial examples on random undercomplete two-layers ReLU neural networks. The term “undercomplete” refers to the fact that their proof only holds when the number of neurons is a vanishing fraction of the ambient dimension. We extend their result to the overcomplete case, where the number of neurons is larger than the dimension (yet also subexponential in the dimension). In fact we prove that a single step of gradient descent suffices. We also show this result for any subexponential width random neural network with smooth activation function.

[Parameterized Knowledge Transfer for Personalized Federated Learning](#)

- Jie Zhang, Song Guo, Xiaosong Ma, Haozhao Wang, Wenchao Xu, Feijie Wu
- abstract@[open-review](#): In recent years, personalized federated learning (pFL) has attracted increasing attention for its potential in dealing with statistical heterogeneity among clients. However, the state-of-the-art pFL methods rely on model parameters aggregation at the server side, which require all models to have the same structure and size, and thus limits the application for more heterogeneous scenarios. To deal with such model constraints, we exploit the potentials of heterogeneous model settings and propose a novel training framework to employ personalized models for different clients. Specifically, we formulate the aggregation procedure in original pFL into a personalized group knowledge transfer training algorithm, namely, KT-pFL, which enables each client to maintain a personalized soft prediction at the server side to guide the others' local training. KT-pFL updates the personalized soft prediction of each client by a linear combination of all local soft predictions using a knowledge coefficient matrix, which can adaptively reinforce the collaboration among clients who own similar data distribution. Furthermore, to quantify the contributions of each client to others' personalized training, the knowledge coefficient matrix is parameterized so that it can be trained simultaneously with the models. The knowledge coefficient matrix and the model parameters are alternatively updated in each round following the gradient descent way. Extensive experiments on various datasets (EMNIST, Fashion_MNIST, CIFAR-10) are conducted under different settings (heterogeneous models and data distributions). It is demonstrated that the proposed framework is the first federated learning paradigm that realizes personalized model training via parameterized group knowledge transfer while achieving significant performance gain comparing with state-of-the-art algorithms.

[Contrastively Disentangled Sequential Variational Autoencoder](#)

- Junwen Bai, Weiran Wang, Carla P. Gomes
- abstract@[open-review](#): Self-supervised disentangled representation learning is a critical task in sequence modeling. The learnt representations contribute to better model interpretability as well as the data generation, and improve the sample efficiency for downstream tasks. We propose a novel sequence representation learning method, named Contrastively Disentangled Sequential Variational Autoencoder (C-DSVAE), to extract and separate the static (time-invariant) and dynamic (time-variant) factors in the latent space. Different from previous sequential variational autoencoder methods, we use a novel evidence lower bound which maximizes the mutual information between the input and the latent factors, while penalizes the mutual information between the static and dynamic factors. We leverage contrastive estimations of the mutual information terms in training, together with simple yet effective augmentation techniques, to introduce additional inductive biases. Our experiments show that C-DSVAE significantly outperforms the previous state-of-the-art methods on multiple metrics.

[Recursive Causal Structure Learning in the Presence of Latent Variables and Selection Bias](#)

- Sina Akbari, Ehsan Mokhtarian, AmirEmad Ghassami, Negar Kiyavash
- abstract@[open-review](#): We consider the problem of learning the causal MAG of a system from observational data in the presence of latent variables and selection bias. Constraint-based methods are one of the main approaches for solving this problem, but the existing methods are either computationally impractical when dealing with large graphs or lacking completeness guarantees. We propose a novel computationally efficient recursive constraint-based method that is sound and complete. The key idea of our approach is that at each iteration a specific type of variable is identified and removed. This allows us to learn the structure efficiently and recursively, as this technique reduces both the number of required conditional independence (CI) tests and the size of the conditioning sets. The former substantially reduces the computational complexity, while the latter results in more reliable CI tests. We provide an upper bound on the number of required CI tests in the worst case. To the best of our knowledge, this is the tightest bound in the literature. We further provide a lower bound on the number of CI tests required by any constraint-based method. The upper bound of our proposed approach and the lower bound at most differ by a factor equal to the number of variables in the worst case. We provide experimental results to compare the proposed approach with the state of the art on both synthetic and real-world structures.

[Generalization Error Rates in Kernel Regression: The Crossover from the Noiseless to Noisy Regime](#)

- Hugo Cui, Bruno Loureiro, Florent Krzakala, Lenka Zdeborová
- abstract@[open-review](#): In this manuscript we consider Kernel Ridge Regression (KRR) under the Gaussian design. Exponents for the decay of the excess generalization error of KRR have been reported in various works under the assumption of power-law decay of eigenvalues of the features co-variance. These decays were, however, provided for sizeably different setups, namely in the noiseless case with constant regularization and in the noisy optimally regularized case. Intermediary settings have been left substantially uncharted. In this work, we unify and extend this line of work, providing characterization of all regimes and excess error decay rates that can be observed in terms of the interplay of noise and regularization. In particular, we show the existence of a transition in the noisy setting between the noiseless exponents to its noisy values as the sample complexity is increased. Finally, we illustrate how this crossover can also be observed on real data sets.

[Learning Gaussian Mixtures with Generalized Linear Models: Precise Asymptotics in High-dimensions](#)

- Bruno Loureiro, Gabriele Sicuro, Cedric Gerbelot, Alessandro Pacco, Florent Krzakala, Lenka Zdeborová
- abstract@[open-review](#): Generalised linear models for multi-class classification problems are one of the fundamental building blocks of modern machine learning tasks. In this manuscript, we characterise the learning of a mixture of K Gaussians with generic means and covariances via empirical risk minimisation (ERM) with any convex loss and regularisation. In particular, we prove exact asymptotics characterising the ERM estimator in high-dimensions, extending several previous results about Gaussian mixture classification in the literature. We exemplify our result in two tasks of interest in statistical learning: a) classification for a mixture with sparse means, where we study the efficiency of ℓ_1 penalty with respect to ℓ_2 ; b) max-margin multi-class classification, where we characterise the phase transition on the existence of the multi-class logistic maximum likelihood estimator for $K > 2$. Finally, we discuss how our theory can be applied beyond the scope of synthetic data, showing that in different cases Gaussian mixtures capture closely the learning curve of classification tasks in real data sets.

Spectral embedding for dynamic networks with stability guarantees

- Ian Gallagher, Andrew Jones, Patrick Rubin-Delanchy
- abstract@[open-review](#): We consider the problem of embedding a dynamic network, to obtain time-evolving vector representations of each node, which can then be used to describe changes in behaviour of individual nodes, communities, or the entire graph. Given this open-ended remit, we argue that two types of stability in the spatio-temporal positioning of nodes are desirable: to assign the same position, up to noise, to nodes behaving similarly at a given time (cross-sectional stability) and a constant position, up to noise, to a single node behaving similarly across different times (longitudinal stability). Similarity in behaviour is defined formally using notions of exchangeability under a dynamic latent position network model. By showing how this model can be recast as a multilayer random dot product graph, we demonstrate that unfolded adjacency spectral embedding satisfies both stability conditions. We also show how two alternative methods, omnibus and independent spectral embedding, alternately lack one or the other form of stability.

Infinite Time Horizon Safety of Bayesian Neural Networks

- Mathias Lechner, Đorđe Žikelić, Krishnendu Chatterjee, Thomas Henzinger
- abstract@[open-review](#): Bayesian neural networks (BNNs) place distributions over the weights of a neural network to model uncertainty in the data and the network's prediction. We consider the problem of verifying safety when running a Bayesian neural network policy in a feedback loop with infinite time horizon systems. Compared to the existing sampling-based approaches, which are inapplicable to the infinite time horizon setting, we train a separate deterministic neural network that serves as an infinite time horizon safety certificate. In particular, we show that the certificate network guarantees the safety of the system over a subset of the BNN weight posterior's support. Our method first computes a safe weight set and then alters the BNN's weight posterior to reject samples outside this set. Moreover, we show how to extend our approach to a safe-exploration reinforcement learning setting, in order to avoid unsafe trajectories during the training of the policy. We evaluate our approach on a series of reinforcement learning benchmarks, including non-Lyapunovian safety specifications.

Towards understanding retrosynthesis by energy-based models

- Ruoxi Sun, Hanjun Dai, Li Li, Steven Kearnes, Bo Dai
- abstract@[open-review](#): Retrosynthesis is the process of identifying a set of reactants to synthesize a target molecule. It is of vital importance to material design and drug discovery. Existing machine learning approaches based on language models and graph neural networks have achieved encouraging results. However, the inner connections of these models are rarely discussed, and rigorous evaluations of these models are largely in need. In this paper, we propose a framework that unifies sequence- and graph-based methods as energy-based models (EBMs) with different energy functions. This unified view establishes connections and reveals the differences between models, thereby enhancing our understanding of model design. We also provide a comprehensive assessment of performance to the community. Moreover, we present a novel dual variant within the framework that performs consistent training to induce the agreement between forward- and backward-prediction. This model improves the state-of-the-art of template-free methods with or without reaction types.

List-Decodable Mean Estimation in Nearly-PCA Time

- Ilias Diakonikolas, Daniel Kane, Daniel Kongsgaard, Jerry Li, Kevin Tian
- abstract@[open-review](#): Robust statistics has traditionally focused on designing estimators tolerant to a minority of contaminated data. {\em List-decodable learning}~\cite{CharikarSV17} studies the more challenging regime where only a minority $\frac{1}{k}$ fraction of the dataset, $k \geq 2$, is drawn from the distribution of interest, and no assumptions are made on the remaining data. We study the fundamental task of list-decodable mean estimation in high dimensions. Our main result is a new algorithm for bounded covariance distributions with optimal sample complexity and near-optimal error guarantee, running in {\em nearly-PCA time}. Assuming the ground truth distribution on \mathbb{R}^d has identity-bounded covariance, our algorithm outputs $O(k)$ candidate means, one of which is within distance $O(\sqrt{k \log k})$ from the truth. Our algorithm runs in time $\tilde{O}(ndk)$, where n is the dataset size. This runtime nearly matches the cost of performing k -PCA on the data, a natural bottleneck of known algorithms for (very) special cases of our problem, such as clustering well-separated mixtures. Prior to our work, the fastest runtimes were $\tilde{O}(n^2 d k^2)$ ~\cite{DiakonikolasKK20}, and $\tilde{O}(nd k^C)$ ~\cite{CherapanamjeriMY20} for an unspecified constant $C \geq 6$. Our approach builds on a novel soft downweighting method we term SIFT, arguably the simplest known polynomial-time mean estimator in the list-decodable setting. To develop our fast algorithms, we boost the computational cost of SIFT via a careful ``win-win-win'' analysis of an approximate Ky Fan matrix multiplicative weights procedure we develop, which may be of independent interest.

Distributed Zero-Order Optimization under Adversarial Noise

- Arya Akhavan, Massimiliano Pontil, Alexandre Tsybakov
- abstract@[open-review](#): We study the problem of distributed zero-order optimization for a class of strongly convex functions. They are formed by the average of local objectives, associated to different nodes in a prescribed network. We propose a distributed zero-order projected gradient descent algorithm to solve the problem. Exchange of information within the network is permitted only between neighbouring nodes. An important feature of our procedure is that it can query only function values, subject to a general noise model, that does not require zero mean or independent errors. We derive upper bounds for the average cumulative regret and optimization error of the algorithm which highlight the role played by a network connectivity parameter, the number of variables, the noise level, the strong convexity parameter, and smoothness properties of the local objectives. The bounds indicate some key improvements of our method over the state-of-the-art, both in the distributed and standard zero-order optimization settings.

Reliable Estimation of KL Divergence using a Discriminator in Reproducing Kernel Hilbert Space

- Sandesh Ghimire, Aria Masoomi, Jennifer Dy
- abstract@[open-review](#): Estimating Kullback–Leibler (KL) divergence from samples of two distributions is essential in many machine learning problems. Variational methods using neural network discriminator have been proposed to achieve this task in a scalable manner. However, we noticed that most of these methods using neural network discriminators suffer from high fluctuations (variance) in estimates and instability in training. In this paper, we look at this issue from statistical learning theory and function space complexity perspective to understand why this happens and how to solve it. We argue that the cause of these pathologies is lack of control over the complexity of the neural network discriminator function and could be mitigated by controlling it. To achieve this objective, we 1) present a novel construction of the discriminator in the Reproducing Kernel Hilbert Space (RKHS), 2) theoretically relate the error probability bound of the KL estimates to the complexity of the discriminator in the RKHS space, 3) present a scalable way to control the complexity (RKHS norm) of the discriminator for a reliable estimation of KL divergence, and 4) prove the consistency of the proposed estimator. In three different applications of KL divergence -- estimation of KL, estimation of mutual information and Variational Bayes -- we show that by controlling the complexity as developed in the theory, we are able to reduce the variance of KL estimates and stabilize the training.

Latent Matters: Learning Deep State-Space Models

- Alexej Klushyn, Richard Kurle, Maximilian Soelch, Botond Cseke, Patrick van der Smagt
- abstract@[open-review](#): Deep state-space models (DSSMs) enable temporal predictions by learning the underlying dynamics of observed sequence data. They are often trained by maximising the evidence lower bound. However, as we show, this does not ensure the model actually learns the underlying

dynamics. We therefore propose a constrained optimisation framework as a general approach for training DSSMs. Building upon this, we introduce the extended Kalman VAE (EKVAE), which combines amortised variational inference with classic Bayesian filtering/smoothing to model dynamics more accurately than RNN-based DSSMs. Our results show that the constrained optimisation framework significantly improves system identification and prediction accuracy on the example of established state-of-the-art DSSMs. The EKVAE outperforms previous models w.r.t. prediction accuracy, achieves remarkable results in identifying dynamical systems, and can furthermore successfully learn state-space representations where static and dynamic features are disentangled.

[On the Estimation Bias in Double Q-Learning](#)

- Zhizhou Ren, Guangxiang Zhu, Hao Hu, Beining Han, Jianglun Chen, Chongjie Zhang
- abstract@[open-review](#): Double Q-learning is a classical method for reducing overestimation bias, which is caused by taking maximum estimated values in the Bellman operation. Its variants in the deep Q-learning paradigm have shown great promise in producing reliable value prediction and improving learning performance. However, as shown by prior work, double Q-learning is not fully unbiased and suffers from underestimation bias. In this paper, we show that such underestimation bias may lead to multiple non-optimal fixed points under an approximate Bellman operator. To address the concerns of converging to non-optimal stationary solutions, we propose a simple but effective approach as a partial fix for the underestimation bias in double Q-learning. This approach leverages an approximate dynamic programming to bound the target value. We extensively evaluate our proposed method in the Atari benchmark tasks and demonstrate its significant improvement over baseline algorithms.

[Mitigating Forgetting in Online Continual Learning with Neuron Calibration](#)

- Haiyan Yin, peng yang, Ping Li
- abstract@[open-review](#): Inspired by human intelligence, the research on online continual learning aims to push the limits of the machine learning models to constantly learn from sequentially encountered tasks, with the data from each task being observed in an online fashion. Though recent studies have achieved remarkable progress in improving the online continual learning performance empowered by the deep neural networks-based models, many of today's approaches still suffer a lot from catastrophic forgetting, a persistent challenge for continual learning. In this paper, we present a novel method which attempts to mitigate catastrophic forgetting in online continual learning from a new perspective, i.e., neuron calibration. In particular, we model the neurons in the deep neural networks-based models as calibrated units under a general formulation. Then we formalize a learning framework to effectively train the calibrated model, where neuron calibration could give ubiquitous benefit to balance the stability and plasticity of online continual learning algorithms through influencing both their forward inference path and backward optimization path. Our proposed formulation for neuron calibration is lightweight and applicable to general feed-forward neural networks-based models. We perform extensive experiments to evaluate our method on four benchmark continual learning datasets. The results show that neuron calibration plays a vital role in improving online continual learning performance and our method could substantially improve the state-of-the-art performance on all~the~evaluated~datasets.

[Escaping Saddle Points with Compressed SGD](#)

- Dmitrii Avdiukhin, Grigory Yaroslavtsev
- abstract@[open-review](#): Stochastic gradient descent (SGD) is a prevalent optimization technique for large-scale distributed machine learning. While SGD computation can be efficiently divided between multiple machines, communication typically becomes a bottleneck in the distributed setting. Gradient compression methods can be used to alleviate this problem, and a recent line of work shows that SGD augmented with gradient compression converges to an $\sqrt{\epsilon}$ -first-order stationary point. In this paper we extend these results to convergence to an $\sqrt{\epsilon}$ -second-order stationary point ($\sqrt{\epsilon}$ -SOSP), which is to the best of our knowledge the first result of this type. In addition, we show that, when the stochastic gradient is not Lipschitz, compressed SGD with RandomK compressor converges to an $\sqrt{\epsilon}$ -SOSP with the same number of iterations as uncompressed SGD [Jin et al.,2021] (JACM), while improving the total communication by a factor of $\tilde{O}(\sqrt{d}\sqrt{\epsilon})$, where d is the dimension of the optimization problem. We present additional results for the cases when the compressor is arbitrary and when the stochastic gradient is Lipschitz.

[Non-Gaussian Gaussian Processes for Few-Shot Regression](#)

- Marcin Sendera, Jacek Tabor, Aleksandra Nowak, Andrzej Bedychaj, Massimiliano Patacchiola, Tomasz Trzcinski, Przemysław Spurek, Maciej Zieba
- abstract@[open-review](#): Gaussian Processes (GPs) have been widely used in machine learning to model distributions over functions, with applications including multi-modal regression, time-series prediction, and few-shot learning. GPs are particularly useful in the last application since they rely on Normal distributions and enable closed-form computation of the posterior probability function. Unfortunately, because the resulting posterior is not flexible enough to capture complex distributions, GPs assume high similarity between subsequent tasks - a requirement rarely met in real-world conditions. In this work, we address this limitation by leveraging the flexibility of Normalizing Flows to modulate the posterior predictive distribution of the GP. This makes the GP posterior locally non-Gaussian, therefore we name our method Non-Gaussian Gaussian Processes (NGGPs). More precisely, we propose an invertible ODE-based mapping that operates on each component of the random variable vectors and shares the parameters across all of them. We empirically tested the flexibility of NGGPs on various few-shot learning regression datasets, showing that the mapping can incorporate context embedding information to model different noise levels for periodic functions. As a result, our method shares the structure of the problem between subsequent tasks, but the contextualization allows for adaptation to dissimilarities. NGGPs outperform the competing state-of-the-art approaches on a diversified set of benchmarks and applications.

[Believe What You See: Implicit Constraint Approach for Offline Multi-Agent Reinforcement Learning](#)

- Yiqin Yang, Xiaoteng Ma, Chenghao Li, Zewu Zheng, Qiyuan Zhang, Gao Huang, Jun Yang, Qianchuan Zhao
- abstract@[open-review](#): Learning from datasets without interaction with environments (Offline Learning) is an essential step to apply Reinforcement Learning (RL) algorithms in real-world scenarios. However, compared with the single-agent counterpart, offline multi-agent RL introduces more agents with the larger state and action space, which is more challenging but attracts little attention. We demonstrate current offline RL algorithms are ineffective in multi-agent systems due to the accumulated extrapolation error. In this paper, we propose a novel offline RL algorithm, named Implicit Constraint Q-learning (ICQ), which effectively alleviates the extrapolation error by only trusting the state-action pairs given in the dataset for value estimation. Moreover, we extend ICQ to multi-agent tasks by decomposing the joint-policy under the implicit constraint. Experimental results demonstrate that the extrapolation error is successfully controlled within a reasonable range and insensitive to the number of agents. We further show that ICQ achieves the state-of-the-art performance in the challenging multi-agent offline tasks (StarCraft II). Our code is public online at <https://github.com/YiqinYang/ICQ>.

[Online Learning in Periodic Zero-Sum Games](#)

- Tanner Fiez, Ryann Sim, Stratis Skoulakis, Georgios Piliouras, Lillian Ratliff
- abstract@[open-review](#): A seminal result in game theory is von Neumann's minmax theorem, which states that zero-sum games admit an essentially unique equilibrium solution. Classical learning results build on this theorem to show that online no-regret dynamics converge to an equilibrium in a time-average sense in zero-sum games. In the past several years, a key research direction has focused on characterizing the transient behavior of such dynamics. General results in this direction show that broad classes of online learning dynamics are cyclic, and formally Poincaré recurrent, in zero-sum games. We analyze the robustness of these online learning behaviors in the case of periodic zero-sum games with a time-invariant equilibrium. This model generalizes the usual repeated game formulation while also being a realistic and natural model of a repeated competition between players that depends on exogenous environmental variations such as time-of-day effects, week-to-week trends, and seasonality. Interestingly, time-average convergence may fail

even in the simplest such settings, in spite of the equilibrium being fixed. In contrast, using novel analysis methods, we show that Poincaré recurrence provably generalizes despite the complex, non-autonomous nature of these dynamical systems.

[K-Net: Towards Unified Image Segmentation](#)

- Wenwei Zhang, Jiangmiao Pang, Kai Chen, Chen Change Loy
- abstract@[open-review](#): Semantic, instance, and panoptic segmentations have been addressed using different and specialized frameworks despite their underlying connections. This paper presents a unified, simple, and effective framework for these essentially similar tasks. The framework, named K-Net, segments both instances and semantic categories consistently by a group of learnable kernels, where each kernel is responsible for generating a mask for either a potential instance or a stuff class. To remedy the difficulties of distinguishing various instances, we propose a kernel update strategy that enables each kernel dynamic and conditional on its meaningful group in the input image. K-Net can be trained in an end-to-end manner with bipartite matching, and its training and inference are naturally NMS-free and box-free. Without bells and whistles, K-Net surpasses all previous published state-of-the-art single-model results of panoptic segmentation on MS COCO test-dev split and semantic segmentation on ADE20K val split with 55.2% PQ and 54.3% mIoU, respectively. Its instance segmentation performance is also on par with Cascade Mask R-CNN on MS COCO with 60%-90% faster inference speeds. Code and models will be released at <https://github.com/ZwwWayne/K-Net/>.

[Pareto-Optimal Learning-Augmented Algorithms for Online Conversion Problems](#)

- Bo Sun, Russell Lee, Mohammad Hajiesmaili, Adam Wierman, Danny Tsang
- abstract@[open-review](#): This paper leverages machine-learned predictions to design competitive algorithms for online conversion problems with the goal of improving the competitive ratio when predictions are accurate (i.e., consistency), while also guaranteeing a worst-case competitive ratio regardless of the prediction quality (i.e., robustness). We unify the algorithmic design of both integral and fractional conversion problems, which are also known as the 1-max-search and one-way trading problems, into a class of online threshold-based algorithms (OTA). By incorporating predictions into design of OTA, we achieve the Pareto-optimal trade-off of consistency and robustness, i.e., no online algorithm can achieve a better consistency guarantee given for a robustness guarantee. We demonstrate the performance of OTA using numerical experiments on Bitcoin conversion.

[Dynaboard: An Evaluation-As-A-Service Platform for Holistic Next-Generation Benchmarking](#)

- Zhiyi Ma, Kawin Ethayarajh, Tristan Thrush, Somya Jain, Ledell Wu, Robin Jia, Christopher Potts, Adina Williams, Douwe Kiela
- abstract@[open-review](#): We introduce Dynaboard, an evaluation-as-a-service framework for hosting benchmarks and conducting holistic model comparison, integrated with the Dynabench platform. Our platform evaluates NLP models directly instead of relying on self-reported metrics or predictions on a single dataset. Under this paradigm, models are submitted to be evaluated in the cloud, circumventing the issues of reproducibility, accessibility, and backwards compatibility that often hinder benchmarking in NLP. This allows users to interact with uploaded models in real time to assess their quality, and permits the collection of additional metrics such as memory use, throughput, and robustness, which -- despite their importance to practitioners -- have traditionally been absent from leaderboards. On each task, models are ranked according to the Dynascore, a novel utility-based aggregation of these statistics, which users can customize to better reflect their preferences, placing more/less weight on a particular axis of evaluation or dataset. As state-of-the-art NLP models push the limits of traditional benchmarks, Dynaboard offers a standardized solution for a more diverse and comprehensive evaluation of model quality.

[NTopo: Mesh-free Topology Optimization using Implicit Neural Representations](#)

- Jonas Zehnder, Yue Li, Stelian Coros, Bernhard Thomaszewski
- abstract@[open-review](#): Recent advances in implicit neural representations show great promise when it comes to generating numerical solutions to partial differential equations. Compared to conventional alternatives, such representations employ parameterized neural networks to define, in a mesh-free manner, signals that are highly-detailed, continuous, and fully differentiable. In this work, we present a novel machine learning approach for topology optimization---an important class of inverse problems with high-dimensional parameter spaces and highly nonlinear objective landscapes. To effectively leverage neural representations in the context of mesh-free topology optimization, we use multilayer perceptrons to parameterize both density and displacement fields. Our experiments indicate that our method is highly competitive for minimizing structural compliance objectives, and it enables self-supervised learning of continuous solution spaces for topology optimization problems.

[Generalization Bounds for \(Wasserstein\) Robust Optimization](#)

- Yang An, Rui Gao
- abstract@[open-review](#): (Distributionally) robust optimization has gained momentum in machine learning community recently, due to its promising applications in developing generalizable learning paradigms. In this paper, we derive generalization bounds for robust optimization and Wasserstein robust optimization for Lipschitz and piecewise Hölder smooth loss functions under both stochastic and adversarial setting, assuming that the underlying data distribution satisfies transportation-information inequalities. The proofs are built on new generalization bounds for variation regularization (such as Lipschitz or gradient regularization) and its connection with robustness.

[Faster Matchings via Learned Duals](#)

- Michael Dinitz, Sungjin Im, Thomas Lavastida, Benjamin Moseley, Sergei Vassilvitskii
- abstract@[open-review](#): A recent line of research investigates how algorithms can be augmented with machine-learned predictions to overcome worst case lower bounds. This area has revealed interesting algorithmic insights into problems, with particular success in the design of competitive online algorithms. However, the question of improving algorithm running times with predictions has largely been unexplored. We take a first step in this direction by combining the idea of machine-learned predictions with the idea of ``warm-starting'' primal-dual algorithms. We consider one of the most important primitives in combinatorial optimization: weighted bipartite matching and its generalization to \$b\$-matching. We identify three key challenges when using learned dual variables in a primal-dual algorithm. First, predicted duals may be infeasible, so we give an algorithm that efficiently maps predicted infeasible duals to nearby feasible solutions. Second, once the duals are feasible, they may not be optimal, so we show that they can be used to quickly find an optimal solution. Finally, such predictions are useful only if they can be learned, so we show that the problem of learning duals for matching has low sample complexity. We validate our theoretical findings through experiments on both real and synthetic data. As a result we give a rigorous, practical, and empirically effective method to compute bipartite matchings.

[Online learning in MDPs with linear function approximation and bandit feedback.](#)

- Gergely Neu, Julia Olkhovskaya
- abstract@[open-review](#): We consider the problem of online learning in an episodic Markov decision process, where the reward function is allowed to change between episodes in an adversarial manner and the learner only observes the rewards associated with its actions. We assume that rewards and the transition function can be represented as linear functions in terms of a known low-dimensional feature map, which allows us to consider the setting where the state space is arbitrarily large. We also assume that the learner has a perfect knowledge of the MDP dynamics. Our main contribution is developing an

algorithm whose expected regret after T episodes is bounded by $\widetilde{O}(\sqrt{dHT})$, where H is the number of steps in each episode and d is the dimensionality of the feature map.

[Learning Collaborative Policies to Solve NP-hard Routing Problems](#)

- Minsu Kim, Jinkyoo Park, joungho kim
- abstract@[open-review](#): Recently, deep reinforcement learning (DRL) frameworks have shown potential for solving NP-hard routing problems such as the traveling salesman problem (TSP) without problem-specific expert knowledge. Although DRL can be used to solve complex problems, DRL frameworks still struggle to compete with state-of-the-art heuristics showing a substantial performance gap. This paper proposes a novel hierarchical problem-solving strategy, termed learning collaborative policies (LCP), which can effectively find the near-optimum solution using two iterative DRL policies: the seeder and reviser. The seeder generates as diversified candidate solutions as possible (seeds) while being dedicated to exploring over the full combinatorial action space (i.e., sequence of assignment action). To this end, we train the seeder's policy using a simple yet effective entropy regularization reward to encourage the seeder to find diverse solutions. On the other hand, the reviser modifies each candidate solution generated by the seeder; it partitions the full trajectory into sub-tours and simultaneously revises each sub-tour to minimize its traveling distance. Thus, the reviser is trained to improve the candidate solution's quality, focusing on the reduced solution space (which is beneficial for exploitation). Extensive experiments demonstrate that the proposed two-policies collaboration scheme improves over single-policy DRL framework on various NP-hard routing problems, including TSP, prize collecting TSP (PCTSP), and capacitated vehicle routing problem (CVRP).

[Efficient Mirror Descent Ascent Methods for Nonsmooth Minimax Problems](#)

- Feihu Huang, Xidong Wu, Heng Huang
- abstract@[open-review](#): In the paper, we propose a class of efficient mirror descent ascent methods to solve the nonsmooth nonconvex-strongly-concave minimax problems by using dynamic mirror functions, and introduce a convergence analysis framework to conduct rigorous theoretical analysis for our mirror descent ascent methods. For our stochastic algorithms, we first prove that the mini-batch stochastic mirror descent ascent (SMDA) method obtains a gradient complexity of $O(\kappa^3\epsilon^{-4})$ for finding an ϵ -stationary point, where κ denotes the condition number. Further, we propose an accelerated stochastic mirror descent ascent (VR-SMDA) method based on the variance reduced technique. We prove that our VR-SMDA method achieves a lower gradient complexity of $O(\kappa^3\epsilon^{-3})$. For our deterministic algorithm, we prove that our deterministic mirror descent ascent (MDA) achieves a lower gradient complexity of $O(\sqrt{\kappa}\epsilon^{-2})$ under mild conditions, which matches the best known complexity in solving smooth nonconvex-strongly-concave minimax optimization. We conduct the experiments on fair classifier and robust neural network training tasks to demonstrate the efficiency of our new algorithms.

[CO-PILOT: Collaborative Planning and reInforcement Learning On sub-Task curriculum](#)

- Shuang Ao, Tianyi Zhou, Guodong Long, Qinghua Lu, Liming Zhu, Jing Jiang
- abstract@[open-review](#): Goal-conditioned reinforcement learning (RL) usually suffers from sparse reward and inefficient exploration in long-horizon tasks. Planning can find the shortest path to a distant goal that provides dense reward/guidance but is inaccurate without a precise environment model. We show that RL and planning can collaboratively learn from each other to overcome their own drawbacks. In "CO-PILOT", a learnable path-planner and an RL agent produce dense feedback to train each other on a curriculum of tree-structured sub-tasks. Firstly, the planner recursively decomposes a long-horizon task to a tree of sub-tasks in a top-down manner, whose layers construct coarse-to-fine sub-task sequences as plans to complete the original task. The planning policy is trained to minimize the RL agent's cost of completing the sequence in each layer from top to bottom layers, which gradually increases the sub-tasks and thus forms an easy-to-hard curriculum for the planner. Next, a bottom-up traversal of the tree trains the RL agent from easier sub-tasks with denser rewards on bottom layers to harder ones on top layers and collects its cost on each sub-task train the planner in the next episode. CO-PILOT repeats this mutual training for multiple episodes before switching to a new task, so the RL agent and planner are fully optimized to facilitate each other's training. We compare CO-PILOT with RL (SAC, HER, PPO), planning (RRT*, NEXT, SGT), and their combination (SoRB) on navigation and continuous control tasks. CO-PILOT significantly improves the success rate and sample efficiency.

[Modality-Agnostic Topology Aware Localization](#)

- Farhad Ghazvinian Zanjani, Ilia Karmanov, Hanno Ackermann, Daniel Dijkman, Simone Merlin, Max Welling, Fatih Porikli
- abstract@[open-review](#): This work presents a data-driven approach for the indoor localization of an observer on a 2D topological map of the environment. State-of-the-art techniques may yield accurate estimates only when they are tailor-made for a specific data modality like camera-based system that prevents their applicability to broader domains. Here, we establish a modality-agnostic framework (called OT-Isomap) and formulate the localization problem in the context of parametric manifold learning while leveraging optimal transportation. This framework allows jointly learning a low-dimensional embedding as well as correspondences with a topological map. We examine the generalizability of the proposed algorithm by applying it to data from diverse modalities such as image sequences and radio frequency signals. The experimental results demonstrate decimeter-level accuracy for localization using different sensory inputs.

[Scalable Quasi-Bayesian Inference for Instrumental Variable Regression](#)

- Ziyu Wang, Yuhao Zhou, Tongzheng Ren, Jun Zhu
- abstract@[open-review](#): Recent years have witnessed an upsurge of interest in employing flexible machine learning models for instrumental variable (IV) regression, but the development of uncertainty quantification methodology is still lacking. In this work we present a scalable quasi-Bayesian procedure for IV regression, building upon the recently developed kernelized IV models. Contrary to Bayesian modeling for IV, our approach does not require additional assumptions on the data generating process, and leads to a scalable approximate inference algorithm with time cost comparable to the corresponding point estimation methods. Our algorithm can be further extended to work with neural network models. We analyze the theoretical properties of the proposed quasi-posterior, and demonstrate through empirical evaluation the competitive performance of our method.

[Kernel Identification Through Transformers](#)

- Fergus Simpson, Ian Davies, Vidhi Lalchand, Alessandro Vullo, Nicolas Durrande, Carl Edward Rasmussen
- abstract@[open-review](#): Kernel selection plays a central role in determining the performance of Gaussian Process (GP) models, as the chosen kernel determines both the inductive biases and prior support of functions under the GP prior. This work addresses the challenge of constructing custom kernel functions for high-dimensional GP regression models. Drawing inspiration from recent progress in deep learning, we introduce a novel approach named KITT: Kernel Identification Through Transformers. KITT exploits a transformer-based architecture to generate kernel recommendations in under 0.1 seconds, which is several orders of magnitude faster than conventional kernel search algorithms. We train our model using synthetic data generated from priors over a vocabulary of known kernels. By exploiting the nature of the self-attention mechanism, KITT is able to process datasets with inputs of arbitrary dimension. We demonstrate that kernels chosen by KITT yield strong performance over a diverse collection of regression benchmarks.

[Curriculum Design for Teaching via Demonstrations: Theory and Applications](#)

- Gaurav Yengeri, Rati Devidze, Parameswaran Kamalaruban, Adish Singla

- abstract@[open-review](#): We consider the problem of teaching via demonstrations in sequential decision-making settings. In particular, we study how to design a personalized curriculum over demonstrations to speed up the learner's convergence. We provide a unified curriculum strategy for two popular learner models: Maximum Causal Entropy Inverse Reinforcement Learning (MaxEnt-IRL) and Cross-Entropy Behavioral Cloning (CrossEnt-BC). Our unified strategy induces a ranking over demonstrations based on a notion of difficulty scores computed w.r.t. the teacher's optimal policy and the learner's current policy. Compared to the state of the art, our strategy doesn't require access to the learner's internal dynamics and still enjoys similar convergence guarantees under mild technical conditions. Furthermore, we adapt our curriculum strategy to the setting where no teacher agent is present using task-specific difficulty scores. Experiments on a synthetic car driving environment and navigation-based environments demonstrate the effectiveness of our curriculum strategy.

[Revenue maximization via machine learning with noisy data](#)

- Ellen Vitercik, Tom Yan
- abstract@[open-review](#): Increasingly, copious amounts of consumer data are used to learn high-revenue mechanisms via machine learning. Existing research on mechanism design via machine learning assumes that there is a distribution over the buyers' values for the items for sale and that the learning algorithm's input is a training set sampled from this distribution. This setup makes the strong assumption that no noise is introduced during data collection. In order to help place mechanism design via machine learning on firm foundations, we investigate the extent to which this learning process is robust to noise. Optimizing revenue using noisy data is challenging because revenue functions are extremely volatile: an infinitesimal change in the buyers' values can cause a steep drop in revenue. Nonetheless, we provide guarantees when arbitrarily correlated noise is added to the training set; we only require that the noise has bounded magnitude or is sub-Gaussian. We conclude with an application of our guarantees to multi-task mechanism design, where there are multiple distributions over buyers' values and the goal is to learn a high-revenue mechanism per distribution. To our knowledge, we are the first to study mechanism design via machine learning with noisy data as well as multi-task mechanism design.

[Exploiting Data Sparsity in Secure Cross-Platform Social Recommendation](#)

- Jinming Cui, Chaochao Chen, Lingjuan Lyu, Carl Yang, Wang Li
- abstract@[open-review](#): Social recommendation has shown promising improvements over traditional systems since it leverages social correlation data as an additional input. Most existing work assumes that all data are available to the recommendation platform. However, in practice, user-item interaction data (e.g., rating) and user-user social data are usually generated by different platforms, and both of which contain sensitive information. Therefore, "How to perform secure and efficient social recommendation across different platforms, where the data are highly-sparse in nature" remains an important challenge. In this work, we bring secure computation techniques into social recommendation, and propose S3Rec, a sparsity-aware secure cross-platform social recommendation framework. As a result, our model can not only improve the recommendation performance of the rating platform by incorporating the sparse social data on the social platform, but also protect data privacy of both platforms. Moreover, to further improve model training efficiency, we propose two secure sparse matrix multiplication protocols based on homomorphic encryption and private information retrieval. Our experiments on two benchmark datasets demonstrate the effectiveness of S3Rec.

[Parallelizing Thompson Sampling](#)

- Amin Karbasi, Vahab Mirrokni, Mohammad Shadravan
- abstract@[open-review](#): How can we make use of information parallelism in online decision-making problems while efficiently balancing the exploration-exploitation trade-off? In this paper, we introduce a batch Thompson Sampling framework for two canonical online decision-making problems with partial feedback, namely, stochastic multi-arm bandit and linear contextual bandit. Over a time horizon $\$T\$$, our batch Thompson Sampling policy achieves the same (asymptotic) regret bound of a fully sequential one while carrying out only $\$O(\log T)\$$ batch queries. To achieve this exponential reduction, i.e., reducing the number of interactions from $\$T\$$ to $\$O(\log T)\$$, our batch policy dynamically determines the duration of each batch in order to balance the exploration-exploitation trade-off. We also demonstrate experimentally that dynamic batch allocation outperforms natural baselines.

[Dynamic Causal Bayesian Optimization](#)

- Virginia Aglietti, Neil Dhir, Javier González, Theodoros Damoulas
- abstract@[open-review](#): We study the problem of performing a sequence of optimal interventions in a dynamic causal system where both the target variable of interest, and the inputs, evolve over time. This problem arises in a variety of domains including healthcare, operational research and policy design. Our approach, which we call Dynamic Causal Bayesian Optimisation (DCBO), brings together ideas from decision making, causal inference and Gaussian process (GP) emulation. DCBO is useful in scenarios where the causal effects are changing over time. Indeed, at every time step, DCBO identifies a local optimal intervention by integrating both observational and past interventional data collected from the system. We give theoretical results detailing how one can transfer interventional information across time steps and define a dynamic causal GP model which can be used to find optimal interventions in practice. Finally, we demonstrate how DCBO identifies optimal interventions faster than competing approaches in multiple settings and applications.

[Local Differential Privacy for Regret Minimization in Reinforcement Learning](#)

- Evrard Garcelon, Vianney Perchet, Ciara Pike-Burke, Matteo Pirotta
- abstract@[open-review](#): Reinforcement learning algorithms are widely used in domains where it is desirable to provide a personalized service. In these domains it is common that user data contains sensitive information that needs to be protected from third parties. Motivated by this, we study privacy in the context of finite-horizon Markov Decision Processes (MDPs) by requiring information to be obfuscated on the user side. We formulate this notion of privacy for RL by leveraging the local differential privacy (LDP) framework. We establish a lower bound for regret minimization in finite-horizon MDPs with LDP guarantees which shows that guaranteeing privacy has a multiplicative effect on the regret. This result shows that while LDP is an appealing notion of privacy, it makes the learning problem significantly more complex. Finally, we present an optimistic algorithm that simultaneously satisfies $\$varepsilon\$$ -LDP requirements, and achieves $\$sqrt{K}varepsilon\$$ regret in any finite-horizon MDP after $\$K\$$ episodes, matching the lower bound dependency on the number of episodes $\$K\$$.

[Emergent Discrete Communication in Semantic Spaces](#)

- Mycal Tucker, Huao Li, Siddharth Agrawal, Dana Hughes, Katia Sycara, Michael Lewis, Julie A Shah
- abstract@[open-review](#): Neural agents trained in reinforcement learning settings can learn to communicate among themselves via discrete tokens, accomplishing as a team what agents would be unable to do alone. However, the current standard of using one-hot vectors as discrete communication tokens prevents agents from acquiring more desirable aspects of communication such as zero-shot understanding. Inspired by word embedding techniques from natural language processing, we propose neural agent architectures that enables them to communicate via discrete tokens derived from a learned, continuous space. We show in a decision theoretic framework that our technique optimizes communication over a wide range of scenarios, whereas one-hot tokens are only optimal under restrictive assumptions. In self-play experiments, we validate that our trained agents learn to cluster tokens in semantically-meaningful ways, allowing them communicate in noisy environments where other techniques fail. Lastly, we demonstrate both that agents using our method can effectively respond to novel human communication and that humans can understand unlabeled emergent agent communication, outperforming the use of one-hot communication.

Drop, Swap, and Generate: A Self-Supervised Approach for Generating Neural Activity

- Ran Liu, Mehdi Azabou, Max Dabagia, Chi-Heng Lin, Mohammad Gheshlaghi Azar, Keith Hengen, Michal Valko, Eva Dyer
- abstract@[open-review](#): Meaningful and simplified representations of neural activity can yield insights into how and what information is being processed within a neural circuit. However, without labels, finding representations that reveal the link between the brain and behavior can be challenging. Here, we introduce a novel unsupervised approach for learning disentangled representations of neural activity called Swap-VAE. Our approach combines a generative modeling framework with an instance-specific alignment loss that tries to maximize the representational similarity between transformed views of the input (brain state). These transformed (or augmented) views are created by dropping out neurons and jittering samples in time, which intuitively should lead the network to a representation that maintains both temporal consistency and invariance to the specific neurons used to represent the neural state. Through evaluations on both synthetic data and neural recordings from hundreds of neurons in different primate brains, we show that it is possible to build representations that disentangle neural datasets along relevant latent dimensions linked to behavior.

Equivariant Manifold Flows

- Isay Katsman, Aaron Lou, Derek Lim, Qingxuan Jiang, Ser Nam Lim, Christopher M. De Sa
- abstract@[open-review](#): Tractably modelling distributions over manifolds has long been an important goal in the natural sciences. Recent work has focused on developing general machine learning models to learn such distributions. However, for many applications these distributions must respect manifold symmetries—a trait which most previous models disregard. In this paper, we lay the theoretical foundations for learning symmetry-invariant distributions on arbitrary manifolds via equivariant manifold flows. We demonstrate the utility of our approach by learning quantum field theory-motivated invariant $SU(n)$ densities and by correcting meteor impact dataset bias.

Scalable Bayesian GPFA with automatic relevance determination and discrete noise models

- Kristopher Jensen, Ta-Chu Kao, Jasmine Stone, Guillaume Hennequin
- abstract@[open-review](#): Latent variable models are ubiquitous in the exploratory analysis of neural population recordings, where they allow researchers to summarize the activity of large populations of neurons in lower dimensional ‘latent’ spaces. Existing methods can generally be categorized into (i) Bayesian methods that facilitate flexible incorporation of prior knowledge and uncertainty estimation, but which typically do not scale to large datasets; and (ii) highly parameterized methods without explicit priors that scale better but often struggle in the low-data regime. Here, we bridge this gap by developing a fully Bayesian yet scalable version of Gaussian process factor analysis (bGPFA), which models neural data as arising from a set of inferred latent processes with a prior that encourages smoothness over time. Additionally, bGPFA uses automatic relevance determination to infer the dimensionality of neural activity directly from the training data during optimization. To enable the analysis of continuous recordings without trial structure, we introduce a novel variational inference strategy that scales near-linearly in time and also allows for non-Gaussian noise models appropriate for electrophysiological recordings. We apply bGPFA to continuous recordings spanning 30 minutes with over 14 million data points from primate motor and somatosensory cortices during a self-paced reaching task. We show that neural activity progresses from an initial state at target onset to a reach-specific preparatory state well before movement onset. The distance between these initial and preparatory latent states is predictive of reaction times across reaches, suggesting that such preparatory dynamics have behavioral relevance despite the lack of externally imposed delay periods. Additionally, bGPFA discovers latent processes that evolve over slow timescales on the order of several seconds and contain complementary information about reaction time. These timescales are longer than those revealed by methods which focus on individual movement epochs and may reflect fluctuations in e.g. task engagement.

Recurrence along Depth: Deep Convolutional Neural Networks with Recurrent Layer Aggregation

- Jingyu Zhao, Yanwen Fang, Guodong Li
- abstract@[open-review](#): This paper introduces a concept of layer aggregation to describe how information from previous layers can be reused to better extract features at the current layer. While DenseNet is a typical example of the layer aggregation mechanism, its redundancy has been commonly criticized in the literature. This motivates us to propose a very light-weighted module, called recurrent layer aggregation (RLA), by making use of the sequential structure of layers in a deep CNN. Our RLA module is compatible with many mainstream deep CNNs, including ResNets, Xception and MobileNetV2, and its effectiveness is verified by our extensive experiments on image classification, object detection and instance segmentation tasks. Specifically, improvements can be uniformly observed on CIFAR, ImageNet and MS COCO datasets, and the corresponding RLA-Nets can surprisingly boost the performances by 2-3% on the object detection task. This evidences the power of our RLA module in helping main CNNs better learn structural information in images.

Independent Prototype Propagation for Zero-Shot Compositionality

- Frank Ruis, Gertjan Burghouts, Doina Bucur
- abstract@[open-review](#): Humans are good at compositional zero-shot reasoning; someone who has never seen a zebra before could nevertheless recognize one when we tell them it looks like a horse with black and white stripes. Machine learning systems, on the other hand, usually leverage spurious correlations in the training data, and while such correlations can help recognize objects in context, they hurt generalization. To be able to deal with underspecified datasets while still leveraging contextual clues during classification, we propose ProtoProp, a novel prototype propagation graph method. First we learn prototypical representations of objects (e.g., zebra) that are independent w.r.t. their attribute labels (e.g., stripes) and vice versa. Next we propagate the independent prototypes through a compositional graph, to learn compositional prototypes of novel attribute-object combinations that reflect the dependencies of the target distribution. The method does not rely on any external data, such as class hierarchy graphs or pretrained word embeddings. We evaluate our approach on AO-Clevr, a synthetic and strongly visual dataset with clean labels, UT-Zappos, a noisy real-world dataset of fine-grained shoe types, and C-GQA, a large-scale object detection dataset modified for compositional zero-shot learning. We show that in the generalized compositional zero-shot setting we outperform state-of-the-art results, and through ablations we show the importance of each part of the method and their contribution to the final results. The code is available on github.

Universal Graph Convolutional Networks

- Di Jin, Zhizhi Yu, Cuiying Huo, Rui Wang, Xiao Wang, Dongxiao He, Jiawei Han
- abstract@[open-review](#): Graph Convolutional Networks (GCNs), aiming to obtain the representation of a node by aggregating its neighbors, have demonstrated great power in tackling various analytics tasks on graph (network) data. The remarkable performance of GCNs typically relies on the homophily assumption of networks, while such assumption cannot always be satisfied, since the heterophily or randomness are also widespread in real-world. This gives rise to one fundamental question: whether networks with different structural properties should adopt different propagation mechanisms? In this paper, we first conduct an experimental investigation. Surprisingly, we discover that there are actually segmentation rules for the propagation mechanism, i.e., 1-hop, 2-hop and \$k\$-nearest neighbor (\$k\$NN) neighbors are more suitable as neighborhoods of network with complete homophily, complete heterophily and randomness, respectively. However, the real-world networks are complex, and may present diverse structural properties, e.g., the network dominated by homophily may contain a small amount of randomness. So can we reasonably utilize these segmentation rules to design a universal propagation mechanism independent of the network structural assumption? To tackle this challenge, we develop a new universal GCN framework, namely U-GCN. It first introduces a multi-type convolution to extract information from 1-hop, 2-hop and \$k\$NN networks simultaneously, and then designs a discriminative aggregation to sufficiently fuse them aiming to given learning objectives. Extensive experiments demonstrate the superiority of U-GCN over state-of-the-arts. The code and data are available at <https://github.com/jindi-tju>.

Adversarial Feature Desensitization

- Pouya Bashivan, Reza Bayat, Adam Ibrahim, Kartik Ahuja, Mojtaba Faramarzi, Touraj Laleh, Blake Richards, Irina Rish
- abstract@[open-review](#): Neural networks are known to be vulnerable to adversarial attacks -- slight but carefully constructed perturbations of the inputs which can drastically impair the network's performance. Many defense methods have been proposed for improving robustness of deep networks by training them on adversarially perturbed inputs. However, these models often remain vulnerable to new types of attacks not seen during training, and even to slightly stronger versions of previously seen attacks. In this work, we propose a novel approach to adversarial robustness, which builds upon the insights from the domain adaptation field. Our method, called Adversarial Feature Desensitization (AFD), aims at learning features that are invariant towards adversarial perturbations of the inputs. This is achieved through a game where we learn features that are both predictive and robust (insensitive to adversarial attacks), i.e. cannot be used to discriminate between natural and adversarial data. Empirical results on several benchmarks demonstrate the effectiveness of the proposed approach against a wide range of attack types and attack strengths. Our code is available at <https://github.com/BashivanLab/afd>.

Few-Shot Data-Driven Algorithms for Low Rank Approximation

- Piotr Indyk, Tal Wagner, David Woodruff
- abstract@[open-review](#): Recently, data-driven and learning-based algorithms for low rank matrix approximation were shown to outperform classical data-oblivious algorithms by wide margins in terms of accuracy. Those algorithms are based on the optimization of sparse sketching matrices, which lead to large savings in time and memory during testing. However, they require long training times on a large amount of existing data, and rely on access to specialized hardware and software. In this work, we develop new data-driven low rank approximation algorithms with better computational efficiency in the training phase, alleviating these drawbacks. Furthermore, our methods are interpretable: while previous algorithms choose the sketching matrix either at random or by black-box learning, we show that it can be set (or initialized) to clearly interpretable values extracted from the dataset. Our experiments show that our algorithms, either by themselves or in combination with previous methods, achieve significant empirical advantage over previous work, improving training times by up to an order of magnitude toward achieving the same target accuracy.

Neural-PIL: Neural Pre-Integrated Lighting for Reflectance Decomposition

- Mark Boss, Varun Jampani, Raphael Braun, Ce Liu, Jonathan Barron, Hendrik PA Lensch
- abstract@[open-review](#): Decomposing a scene into its shape, reflectance and illumination is a fundamental problem in computer vision and graphics. Neural approaches such as NeRF have achieved remarkable success in view synthesis, but do not explicitly perform decomposition and instead operate exclusively on radiance (the product of reflectance and illumination). Extensions to NeRF, such as NeRD, can perform decomposition but struggle to accurately recover detailed illumination, thereby significantly limiting realism. We propose a novel reflectance decomposition network that can estimate shape, BRDF, and per-image illumination given a set of object images captured under varying illumination. Our key technique is a novel illumination integration network called Neural-PIL that replaces a costly illumination integral operation in the rendering with a simple network query. In addition, we also learn deep low-dimensional priors on BRDF and illumination representations using novel smooth manifold auto-encoders. Our decompositions can result in considerably better BRDF and light estimates enabling more accurate novel view-synthesis and relighting compared to prior art. Project page: <https://markboss.me/publication/2021-neural-pil/>

Asymptotics of the Bootstrap via Stability with Applications to Inference with Model Selection

- Morgane Austern, Vasilis Syrgkanis
- abstract@[open-review](#): One of the most commonly used methods for forming confidence intervals is the empirical bootstrap, which is especially expedient when the limiting distribution of the estimator is unknown. However, despite its ubiquitous role in machine learning, its theoretical properties are still not well understood. Recent developments in probability have provided new tools to study the bootstrap method. However, they have been applied only to specific applications and contexts, and it is unclear whether these techniques are applicable to the understanding of the consistency of the bootstrap in machine learning pipelines. In this paper, we derive general stability conditions under which the empirical bootstrap estimator is consistent and quantify the speed of convergence. Moreover, we propose alternative ways to use the bootstrap method to build confidence intervals with coverage guarantees. Finally, we illustrate the generality and tightness of our results by examples of interest for machine learning including for two-sample kernel tests after kernel selection and the empirical risk of stacked estimators.

Dynamic influence maximization

- Binghui Peng
- abstract@[open-review](#): We initiate a systematic study on {\em dynamic influence maximization} (DIM). In the DIM problem, one maintains a seed set \mathcal{S} of at most k nodes in a dynamically involving social network, with the goal of maximizing the expected influence spread while minimizing the amortized updating cost. We consider two evolution models. In the {\em incremental model}, the social network gets enlarged over time and one only introduces new users and establishes new social links, we design an algorithm that achieves $(1-1/e-\epsilon)$ -approximation to the optimal solution and has $k \cdot \text{poly}(\log n, \epsilon^{-1})$ amortized running time, which matches the state-of-art offline algorithm with only poly-logarithmic overhead. In the fully dynamic model, users join in and leave, influence propagation gets strengthened or weakened in real time, we prove that under the Strong Exponential Time Hypothesis (SETH), no algorithm can achieve $2^{-(\log n)^{1-o(1)}}\text{-approximation}$ unless the amortized running time is $n^{1-o(1)}$. On the technical side, we exploit novel adaptive sampling approaches that reduce DIM to the dynamic MAX-k coverage problem, and design an efficient $(1-1/e-\epsilon)$ -approximation algorithm for it. Our lower bound leverages the recent developed distributed PCP framework.

Risk Monotonicity in Statistical Learning

- Zakaria Mhammedi
- abstract@[open-review](#): Acquisition of data is a difficult task in many applications of machine learning, and it is only natural that one hopes and expects the population risk to decrease (better performance) monotonically with increasing data points. It turns out, somewhat surprisingly, that this is not the case even for the most standard algorithms that minimize the empirical risk. Non-monotonic behavior of the risk and instability in training have manifested and appeared in the popular deep learning paradigm under the description of double descent. These problems highlight the current lack of understanding of learning algorithms and generalization. It is, therefore, crucial to pursue this concern and provide a characterization of such behavior. In this paper, we derive the first consistent and risk-monotonic (in high probability) algorithms for a general statistical learning setting under weak assumptions, consequently answering some questions posed by Viering et. al. 2019 on how to avoid non-monotonic behavior of risk curves. We further show that risk monotonicity need not necessarily come at the price of worse excess risk rates. To achieve this, we derive new empirical Bernstein-like concentration inequalities of independent interest that hold for certain non-i.i.d.-processes such as Martingale Difference Sequences.

Information is Power: Intrinsic Control via Information Capture

- Nicholas Rhinehart, Jenny Wang, Glen Berseth, John Co-Reyes, Danijar Hafner, Chelsea Finn, Sergey Levine
- abstract@[open-review](#): Humans and animals explore their environment and acquire useful skills even in the absence of clear goals, exhibiting intrinsic motivation. The study of intrinsic motivation in artificial agents is concerned with the following question: what is a good general-purpose objective for an

agent? We study this question in dynamic partially-observed environments, and argue that a compact and general learning objective is to minimize the entropy of the agent's state visitation estimated using a latent state-space model. This objective induces an agent to both gather information about its environment, corresponding to reducing uncertainty, and to gain control over its environment, corresponding to reducing the unpredictability of future world states. We instantiate this approach as a deep reinforcement learning agent equipped with a deep variational Bayes filter. We find that our agent learns to discover, represent, and exercise control of dynamic objects in a variety of partially-observed environments sensed with visual observations without extrinsic reward.

[Extracting Deformation-Aware Local Features by Learning to Deform](#)

- Guilherme Potje, Renato Martins, Felipe Chamone, Erickson Nascimento
- abstract@[open-review](#): Despite the advances in extracting local features achieved by handcrafted and learning-based descriptors, they are still limited by the lack of invariance to non-rigid transformations. In this paper, we present a new approach to compute features from still images that are robust to non-rigid deformations to circumvent the problem of matching deformable surfaces and objects. Our deformation-aware local descriptor, named DEAL, leverages a polar sampling and a spatial transformer warping to provide invariance to rotation, scale, and image deformations. We train the model architecture end-to-end by applying isometric non-rigid deformations to objects in a simulated environment as guidance to provide highly discriminative local features. The experiments show that our method outperforms state-of-the-art handcrafted, learning-based image, and RGB-D descriptors in different datasets with both real and realistic synthetic deformable objects in still images. The source code and trained model of the descriptor are publicly available at <https://www.verlab.dcc.ufmg.br/descriptors/neurips2021>.

[Object-Centric Representation Learning with Generative Spatial-Temporal Factorization](#)

- Nanbo Li, Muhammad Ahmed Raza, Wenbin Hu, Zhaole Sun, Robert Fisher
- abstract@[open-review](#): Learning object-centric scene representations is essential for attaining structural understanding and abstraction of complex scenes. Yet, as current approaches for unsupervised object-centric representation learning are built upon either a stationary observer assumption or a static scene assumption, they often: i) suffer single-view spatial ambiguities, or ii) infer incorrectly or inaccurately object representations from dynamic scenes. To address this, we propose Dynamics-aware Multi-Object Network (DyMON), a method that broadens the scope of multi-view object-centric representation learning to dynamic scenes. We train DyMON on multi-view-dynamic-scene data and show that DyMON learns---without supervision---to factorize the entangled effects of observer motions and scene object dynamics from a sequence of observations, and constructs scene object spatial representations suitable for rendering at arbitrary times (querying across time) and from arbitrary viewpoints (querying across space). We also show that the factorized scene representations (w.r.t. objects) support querying about a single object by space and time independently.

[Learning to Simulate Self-driven Particles System with Coordinated Policy Optimization](#)

- Zhenghao Peng, Quanyi Li, Ka Ming Hui, Chunxiao Liu, Bolei Zhou
- abstract@[open-review](#): Self-Driven Particles (SDP) describe a category of multi-agent systems common in everyday life, such as flocking birds and traffic flows. In a SDP system, each agent pursues its own goal and constantly changes its cooperative or competitive behaviors with its nearby agents. Manually designing the controllers for such SDP system is time-consuming, while the resulting emergent behaviors are often not realistic nor generalizable. Thus the realistic simulation of SDP systems remains challenging. Reinforcement learning provides an appealing alternative for automating the development of the controller for SDP. However, previous multi-agent reinforcement learning (MARL) methods define the agents to be teammates or enemies before hand, which fail to capture the essence of SDP where the role of each agent varies to be cooperative or competitive even within one episode. To simulate SDP with MARL, a key challenge is to coordinate agents' behaviors while still maximizing individual objectives. Taking traffic simulation as the testing bed, in this work we develop a novel MARL method called Coordinated Policy Optimization (CoPO), which incorporates social psychology principle to learn neural controller for SDP. Experiments show that the proposed method can achieve superior performance compared to MARL baselines in various metrics. Noticeably the trained vehicles exhibit complex and diverse social behaviors that improve performance and safety of the population as a whole. Demo video and source code are available at: <https://decisionforce.github.io/CoPO/>

[Gradient-based Hyperparameter Optimization Over Long Horizons](#)

- Paul Micaelli, Amos J. Storkey
- abstract@[open-review](#): Gradient-based hyperparameter optimization has earned a widespread popularity in the context of few-shot meta-learning, but remains broadly impractical for tasks with long horizons (many gradient steps), due to memory scaling and gradient degradation issues. A common workaround is to learn hyperparameters online, but this introduces greediness which comes with a significant performance drop. We propose forward-mode differentiation with sharing (FDS), a simple and efficient algorithm which tackles memory scaling issues with forward-mode differentiation, and gradient degradation issues by sharing hyperparameters that are contiguous in time. We provide theoretical guarantees about the noise reduction properties of our algorithm, and demonstrate its efficiency empirically by differentiating through $\sim 10^4$ gradient steps of unrolled optimization. We consider large hyperparameter search ranges on CIFAR-10 where we significantly outperform greedy gradient-based alternatives, while achieving $\times 20$ speedups compared to the state-of-the-art black-box methods.

[Stochastic Bias-Reduced Gradient Methods](#)

- Hilal Asi, Yair Carmon, Arun Jambulapati, Yujia Jin, Aaron Sidford
- abstract@[open-review](#): We develop a new primitive for stochastic optimization: a low-bias, low-cost estimator of the minimizer x_{\star} of any Lipschitz strongly-convex function f . In particular, we use a multilevel Monte-Carlo approach due to Blanchet and Glynn to turn any optimal stochastic gradient method into an estimator of x_{\star} with bias δ , variance $O(\log(1/\delta))$, and an expected sampling cost of $O(\log(1/\delta))$ stochastic gradient evaluations. As an immediate consequence, we obtain cheap and nearly unbiased gradient estimators for the Moreau envelope of any Lipschitz convex function. We demonstrate the potential of our estimator through four applications. First, we develop a method for minimizing the maximum of N functions, improving on recent results and matching a lower bound up to logarithmic factors. Second and third, we recover state-of-the-art rates for projection-efficient and gradient-efficient optimization using simple algorithms with a transparent analysis. Finally, we show that an improved version of our estimator would yield a nearly linear-time, optimal-utility, differentially-private non-smooth stochastic optimization method.

[The Causal-Neural Connection: Expressiveness, Learnability, and Inference](#)

- Kevin Xia, Kai-Zhan Lee, Yoshua Bengio, Elias Bareinboim
- abstract@[open-review](#): One of the central elements of any causal inference is an object called structural causal model (SCM), which represents a collection of mechanisms and exogenous sources of random variation of the system under investigation (Pearl, 2000). An important property of many kinds of neural networks is universal approximability: the ability to approximate any function to arbitrary precision. Given this property, one may be tempted to surmise that a collection of neural nets is capable of learning any SCM by training on data generated by that SCM. In this paper, we show this is not the case by disentangling the notions of expressivity and learnability. Specifically, we show that the causal hierarchy theorem (Thm. 1, Bareinboim et al., 2020), which describes the limits of what can be learned from data, still holds for neural models. For instance, an arbitrarily complex and expressive neural net is unable to predict the effects of interventions given observational data alone. Given this result, we introduce a special type of SCM called a neural causal model (NCM), and formalize a new type of inductive bias to encode structural constraints necessary for performing causal inferences. Building on this new class of models, we focus on solving two canonical tasks found in the literature known as causal identification and estimation.

Leveraging the neural toolbox, we develop an algorithm that is both sufficient and necessary to determine whether a causal effect can be learned from data (i.e., causal identifiability); it then estimates the effect whenever identifiability holds (causal estimation). Simulations corroborate the proposed approach.

[Validation Free and Replication Robust Volume-based Data Valuation](#)

- Xinyi Xu, Zhaoxuan Wu, Chuan Sheng Foo, Bryan Kian Hsiang Low
- abstract@[open-review](#): Data valuation arises as a non-trivial challenge in real-world use cases such as collaborative machine learning, federated learning, trusted data sharing, data marketplaces. The value of data is often associated with the learning performance (e.g., validation accuracy) of a model trained on the data, which introduces a close coupling between data valuation and validation. However, a validation set may not be available in practice and it can be challenging for the data providers to reach an agreement on the choice of the validation set. Another practical issue is that of data replication: Given the value of some data points, a dishonest data provider may replicate these data points to exploit the valuation for a larger reward/payment. We observe that the diversity of the data points is an inherent property of a dataset that is independent of validation. We formalize diversity via the volume of the data matrix (i.e., determinant of its left Gram), which allows us to establish a formal connection between the diversity of data and learning performance without requiring validation. Furthermore, we propose a robust volume measure with a theoretical guarantee on the replication robustness by following the intuition that copying the same data points does not increase the diversity of data. We perform extensive experiments to demonstrate its consistency in valuation and practical advantages over existing baselines and show that our method is model- and task-agnostic and can be flexibly adapted to handle various neural networks.

[Implicit Finite-Horizon Approximation and Efficient Optimal Algorithms for Stochastic Shortest Path](#)

- Liyu Chen, Mehdi Jafarnia-Jahromi, Rahul Jain, Haipeng Luo
- abstract@[open-review](#): We introduce a generic template for developing regret minimization algorithms in the Stochastic Shortest Path (SSP) model, which achieves minimax optimal regret as long as certain properties are ensured. The key of our analysis is a new technique called implicit finite-horizon approximation, which approximates the SSP model by a finite-horizon counterpart only in the analysis without explicit implementation. Using this template, we develop two new algorithms: the first one is model-free (the first in the literature to our knowledge) and minimax optimal under strictly positive costs; the second one is model-based and minimax optimal even with zero-cost state-action pairs, matching the best existing result from [Tarbouriech et al., 2021b]. Importantly, both algorithms admit highly sparse updates, making them computationally more efficient than all existing algorithms. Moreover, both can be made completely parameter-free.

[A Separation Result Between Data-oblivious and Data-aware Poisoning Attacks](#)

- Samuel Deng, Sanjam Garg, Somesh Jha, Saeed Mahloujifar, Mohammad Mahmoody, Abhradeep Guha Thakurta
- abstract@[open-review](#): Poisoning attacks have emerged as a significant security threat to machine learning algorithms. It has been demonstrated that adversaries who make small changes to the training set, such as adding specially crafted data points, can hurt the performance of the output model. Most of these attacks require the full knowledge of training data. This leaves open the possibility of achieving the same attack results using poisoning attacks that do not have the full knowledge of the clean training set. In this work, we initiate a theoretical study of the problem above. Specifically, for the case of feature selection with LASSO, we show that *full information* adversaries (that craft poisoning examples based on the rest of the training data) are provably much more devastating compared to the optimal attacker that is *oblivious* to the training set yet has access to the distribution of the data. Our separation result shows that the two settings of data-aware and data-oblivious are fundamentally different and we cannot hope to achieve the same attack or defense results in these scenarios.

[Deep Learning Through the Lens of Example Difficulty](#)

- Robert Baldock, Hartmut Maennel, Behnam Neyshabur
- abstract@[open-review](#): Existing work on understanding deep learning often employs measures that compress all data-dependent information into a few numbers. In this work, we adopt a perspective based on the role of individual examples. We introduce a measure of the computational difficulty of making a prediction for a given input: the (effective) prediction depth. Our extensive investigation reveals surprising yet simple relationships between the prediction depth of a given input and the model's uncertainty, confidence, accuracy and speed of learning for that data point. We further categorize difficult examples into three interpretable groups, demonstrate how these groups are processed differently inside deep models and showcase how this understanding allows us to improve prediction accuracy. Insights from our study lead to a coherent view of a number of separately reported phenomena in the literature: early layers generalize while later layers memorize; early layers converge faster and networks learn easy data and simple functions first.

[R-Drop: Regularized Dropout for Neural Networks](#)

- xiaobo liang, Lijun Wu, Juntao Li, Yue Wang, Qi Meng, Tao Qin, Wei Chen, Min Zhang, Tie-Yan Liu
- abstract@[open-review](#): Dropout is a powerful and widely used technique to regularize the training of deep neural networks. Though effective and performing well, the randomness introduced by dropout causes unnegligible inconsistency between training and inference. In this paper, we introduce a simple consistency training strategy to regularize dropout, namely R-Drop, which forces the output distributions of different sub models generated by dropout to be consistent with each other. Specifically, for each training sample, R-Drop minimizes the bidirectional KL-divergence between the output distributions of two sub models sampled by dropout. Theoretical analysis reveals that R-Drop reduces the above inconsistency. Experiments on \$\\bf{5}\$ widely used deep learning tasks (\$\\bf{18}\$ datasets in total), including neural machine translation, abstractive summarization, language understanding, language modeling, and image classification, show that R-Drop is universally effective. In particular, it yields substantial improvements when applied to fine-tune large-scale pre-trained models, e.g., ViT, RoBERTa-large, and BART, and achieves state-of-the-art (SOTA) performances with the vanilla Transformer model on WMT14 English\$\\rightarrow\$German translation (\$\\bf{30.91}\$ BLEU) and WMT14 English\$\\rightarrow\$French translation (\$\\bf{43.95}\$ BLEU), even surpassing models trained with extra large-scale data and expert-designed advanced variants of Transformer models. Our code is available at GitHub\\footnote{\$\\{\\url{https://github.com/dropreg/R-Drop}\\}\$}.

[Diversity Enhanced Active Learning with Strictly Proper Scoring Rules](#)

- Wei Tan, Lan Du, Wray Buntine
- abstract@[open-review](#): We study acquisition functions for active learning (AL) for text classification. The Expected Loss Reduction (ELR) method focuses on a Bayesian estimate of the reduction in classification error, recently updated with Mean Objective Cost of Uncertainty (MOCU). We convert the ELR framework to estimate the increase in (strictly proper) scores like log probability or negative mean square error, which we call Bayesian Estimate of Mean Proper Scores (BEMPS). We also prove convergence results borrowing techniques used with MOCU. In order to allow better experimentation with the new acquisition functions, we develop a complementary batch AL algorithm, which encourages diversity in the vector of expected changes in scores for unlabelled data. To allow high performance text classifiers, we combine ensembling and dynamic validation set construction on pretrained language models. Extensive experimental evaluation then explores how these different acquisition functions perform. The results show that the use of mean square error and log probability with BEMPS yields robust acquisition functions, which consistently outperform the others tested.

[SSUL: Semantic Segmentation with Unknown Label for Exemplar-based Class-Incremental Learning](#)

- Sungmin Cha, beomyoung kim, YoungJoon Yoo, Taesup Moon
- abstract@[open-review](#): We consider a class-incremental semantic segmentation (CISS) problem. While some recently proposed algorithms utilized variants of knowledge distillation (KD) technique to tackle the problem, they only partially addressed the key additional challenges in CISS that causes the catastrophic forgetting; \textit{i.e.}, the semantic drift of the background class and multi-label prediction issue. To better address these challenges, we propose a new method, dubbed as SSUL-M (Semantic Segmentation with Unknown Label with Memory), by carefully combining several techniques tailored for semantic segmentation. More specifically, we make three main contributions; (1) modeling \textit{unknown} class within the background class to help learning future classes (help plasticity), (2) \textit{freezing} backbone network and past classifiers with binary cross-entropy loss and pseudo-labeling to overcome catastrophic forgetting (help stability), and (3) utilizing \textit{tiny exemplar memory} for the first time in CISS to improve \textit{both} plasticity and stability. As a result, we show our method achieves significantly better performance than the recent state-of-the-art baselines on the standard benchmark datasets. Furthermore, we justify our contributions with thorough and extensive ablation analyses and discuss different natures of the CISS problem compared to the standard class-incremental learning for classification. The official code is available at <https://github.com/clovaai/SSUL>.

[Lower and Upper Bounds on the Pseudo-Dimension of Tensor Network Models](#)

- Behroush Khavari, Guillaume Rabusseau
- abstract@[open-review](#): Tensor network methods have been a key ingredient of advances in condensed matter physics and have recently sparked interest in the machine learning community for their ability to compactly represent very high-dimensional objects. Tensor network methods can for example be used to efficiently learn linear models in exponentially large feature spaces [Stoudenmire and Schwab, 2016]. In this work, we derive upper and lower bounds on the VC dimension and pseudo-dimension of a large class of tensor network models for classification, regression and completion. Our upper bounds hold for linear models parameterized by arbitrary tensor network structures, and we derive lower bounds for common tensor decomposition models~(CP, Tensor Train, Tensor Ring and Tucker) showing the tightness of our general upper bound. These results are used to derive a generalization bound which can be applied to classification with low rank matrices as well as linear classifiers based on any of the commonly used tensor decomposition models. As a corollary of our results, we obtain a bound on the VC dimension of the matrix product state classifier introduced in [Stoudenmire and Schwab, 2016] as a function of the so-called bond dimension~(i.e. tensor train rank), which answers an open problem listed by Cirac, Garre-Rubio and Pérez-García in [Cirac et al., 2019].

[What Makes Multi-Modal Learning Better than Single \(Provably\)](#)

- Yu Huang, Chenzhuang Du, Zihui Xue, Xuanyao Chen, Hang Zhao, Longbo Huang
- abstract@[open-review](#): The world provides us with data of multiple modalities. Intuitively, models fusing data from different modalities outperform their uni-modal counterparts, since more information is aggregated. Recently, joining the success of deep learning, there is an influential line of work on deep multi-modal learning, which has remarkable empirical results on various applications. However, theoretical justifications in this field are notably lacking. Can multi-modal learning provably perform better than uni-modal? In this paper, we answer this question under a most popular multi-modal fusion framework, which firstly encodes features from different modalities into a common latent space and seamlessly maps the latent representations into the task space. We prove that learning with multiple modalities achieves a smaller population risk than only using its subset of modalities. The main intuition is that the former has a more accurate estimate of the latent space representation. To the best of our knowledge, this is the first theoretical treatment to capture important qualitative phenomena observed in real multi-modal applications from the generalization perspective. Combining with experiment results, we show that multi-modal learning does possess an appealing formal guarantee.

[Quantifying and Improving Transferability in Domain Generalization](#)

- Guojun Zhang, Han Zhao, Yaoliang Yu, Pascal Poupart
- abstract@[open-review](#): Out-of-distribution generalization is one of the key challenges when transferring a model from the lab to the real world. Existing efforts mostly focus on building invariant features among source and target domains. Based on invariant features, a high-performing classifier on source domains could hopefully behave equally well on a target domain. In other words, we hope the invariant features to be \emph{transferable}. However, in practice, there are no perfectly transferable features, and some algorithms seem to learn ``more transferable'' features than others. How can we understand and quantify such \emph{transferability}? In this paper, we formally define transferability that one can quantify and compute in domain generalization. We point out the difference and connection with common discrepancy measures between domains, such as total variation and Wasserstein distance. We then prove that our transferability can be estimated with enough samples and give a new upper bound for the target error based on our transferability. Empirically, we evaluate the transferability of the feature embeddings learned by existing algorithms for domain generalization. Surprisingly, we find that many algorithms are not quite learning transferable features, although few could still survive. In light of this, we propose a new algorithm for learning transferable features and test it over various benchmark datasets, including RotatedMNIST, PACS, Office-Home and WILDS-FMoW. Experimental results show that the proposed algorithm achieves consistent improvement over many state-of-the-art algorithms, corroborating our theoretical findings.

[Beyond Pinball Loss: Quantile Methods for Calibrated Uncertainty Quantification](#)

- Youngseog Chung, Willie Neiswanger, Ian Char, Jeff Schneider
- abstract@[open-review](#): Among the many ways of quantifying uncertainty in a regression setting, specifying the full quantile function is attractive, as quantiles are amenable to interpretation and evaluation. A model that predicts the true conditional quantiles for each input, at all quantile levels, presents a correct and efficient representation of the underlying uncertainty. To achieve this, many current quantile-based methods focus on optimizing the pinball loss. However, this loss restricts the scope of applicable regression models, limits the ability to target many desirable properties (e.g. calibration, sharpness, centered intervals), and may produce poor conditional quantiles. In this work, we develop new quantile methods that address these shortcomings. In particular, we propose methods that can apply to any class of regression model, select an explicit balance between calibration and sharpness, optimize for calibration of centered intervals, and produce more accurate conditional quantiles. We provide a thorough experimental evaluation of our methods, which includes a high dimensional uncertainty quantification task in nuclear fusion.

[Dynamic Inference with Neural Interpreters](#)

- Nasim Rahaman, Muhammad Waleed Gondal, Shruti Joshi, Peter Gehler, Yoshua Bengio, Francesco Locatello, Bernhard Schölkopf
- abstract@[open-review](#): Modern neural network architectures can leverage large amounts of data to generalize well within the training distribution. However, they are less capable of systematic generalization to data drawn from unseen but related distributions, a feat that is hypothesized to require compositional reasoning and reuse of knowledge. In this work, we present Neural Interpreters, an architecture that factorizes inference in a self-attention network as a system of modules, which we call functions. Inputs to the model are routed through a sequence of functions in a way that is end-to-end learned. The proposed architecture can flexibly compose computation along width and depth, and lends itself well to capacity extension after training. To demonstrate the versatility of Neural Interpreters, we evaluate it in two distinct settings: image classification and visual abstract reasoning on Raven Progressive Matrices. In the former, we show that Neural Interpreters perform on par with the vision transformer using fewer parameters, while being transferrable to a new task in a sample efficient manner. In the latter, we find that Neural Interpreters are competitive with respect to the state-of-the-art in terms of systematic generalization.

[Leveraging Recursive Gumbel-Max Trick for Approximate Inference in Combinatorial Spaces](#)

- Kirill Struminsky, Artyom Gadetsky, Denis Rakitin, Danil Karpushkin, Dmitry P. Vetrov

- abstract@[open-review](#): Structured latent variables allow incorporating meaningful prior knowledge into deep learning models. However, learning with such variables remains challenging because of their discrete nature. Nowadays, the standard learning approach is to define a latent variable as a perturbed algorithm output and to use a differentiable surrogate for training. In general, the surrogate puts additional constraints on the model and inevitably leads to biased gradients. To alleviate these shortcomings, we extend the Gumbel-Max trick to define distributions over structured domains. We avoid the differentiable surrogates by leveraging the score function estimators for optimization. In particular, we highlight a family of recursive algorithms with a common feature we call stochastic invariant. The feature allows us to construct reliable gradient estimates and control variates without additional constraints on the model. In our experiments, we consider various structured latent variable models and achieve results competitive with relaxation-based counterparts.

[Hamiltonian Dynamics with Non-Newtonian Momentum for Rapid Sampling](#)

- Greg Ver Steeg, Aram Galstyan
- abstract@[open-review](#): Sampling from an unnormalized probability distribution is a fundamental problem in machine learning with applications including Bayesian modeling, latent factor inference, and energy-based model training. After decades of research, variations of MCMC remain the default approach to sampling despite slow convergence. Auxiliary neural models can learn to speed up MCMC, but the overhead for training the extra model can be prohibitive. We propose a fundamentally different approach to this problem via a new Hamiltonian dynamics with a non-Newtonian momentum. In contrast to MCMC approaches like Hamiltonian Monte Carlo, no stochastic step is required. Instead, the proposed deterministic dynamics in an extended state space exactly sample the target distribution, specified by an energy function, under an assumption of ergodicity. Alternatively, the dynamics can be interpreted as a normalizing flow that samples a specified energy model without training. The proposed Energy Sampling Hamiltonian (ESH) dynamics have a simple form that can be solved with existing ODE solvers, but we derive a specialized solver that exhibits much better performance. ESH dynamics converge faster than their MCMC competitors enabling faster, more stable training of neural network energy models.

[Dynamic Normalization and Relay for Video Action Recognition](#)

- Dongqi Cai, Anbang Yao, Yurong Chen
- abstract@[open-review](#): Convolutional Neural Networks (CNNs) have been the dominant model for video action recognition. Due to the huge memory and compute demand, popular action recognition networks need to be trained with small batch sizes, which makes learning discriminative spatial-temporal representations for videos become a challenging problem. In this paper, we present Dynamic Normalization and Relay (DNR), an improved normalization design, to augment the spatial-temporal representation learning of any deep action recognition model, adapting to small batch size training settings. We observe that state-of-the-art action recognition networks usually apply the same normalization parameters to all video data, and ignore the dependencies of the estimated normalization parameters between neighboring frames (at the same layer) and between neighboring layers (with all frames of a video clip). Inspired by this, DNR introduces two dynamic normalization relay modules to explore the potentials of cross-temporal and cross-layer feature distribution dependencies for estimating accurate layer-wise normalization parameters. These two DNR modules are instantiated as a light-weight recurrent structure conditioned on the current input features, and the normalization parameters estimated from the neighboring frames based features at the same layer or from the whole video clip based features at the preceding layers. We first plug DNR into prevailing 2D CNN backbones and test its performance on public action recognition datasets including Kinetics and Something-Something. Experimental results show that DNR brings large performance improvements to the baselines, achieving over 4.4% absolute margins in top-1 accuracy without training bells and whistles. More experiments on 3D backbones and several latest 2D spatial-temporal networks further validate its effectiveness. Code will be available at <https://github.com/caidonkey/dnr>.

[Robust Visual Reasoning via Language Guided Neural Module Networks](#)

- Arjun Akula, Varun Jampani, Soravit Changpinyo, Song-Chun Zhu
- abstract@[open-review](#): Neural module networks (NMN) are a popular approach for solving multi-modal tasks such as visual question answering (VQA) and visual referring expression recognition (REF). A key limitation in prior implementations of NMN is that the neural modules do not effectively capture the association between the visual input and the relevant neighbourhood context of the textual input. This limits their generalizability. For instance, NMN fail to understand new concepts such as "yellow sphere to the left" even when it is a combination of known concepts from train data: "blue sphere", "yellow cube", and "metallic cube to the left". In this paper, we address this limitation by introducing a language-guided adaptive convolution layer (LG-Conv) into NMN, in which the filter weights of convolutions are explicitly multiplied with a spatially varying language-guided kernel. Our model allows the neural module to adaptively co-attend over potential objects of interest from the visual and textual inputs. Extensive experiments on VQA and REF tasks demonstrate the effectiveness of our approach. Additionally, we propose a new challenging out-of-distribution test split for REF task, which we call C3-Ref+, for explicitly evaluating the NMN's ability to generalize well to adversarial perturbations and unseen combinations of known concepts. Experiments on C3-Ref+ further demonstrate the generalization capabilities of our approach.

[True Few-Shot Learning with Language Models](#)

- Ethan Perez, Douwe Kiela, Kyunghyun Cho
- abstract@[open-review](#): Pretrained language models (LMs) perform well on many tasks even when learning from a few examples, but prior work uses many held-out examples to tune various aspects of learning, such as hyperparameters, training objectives, and natural language templates ("prompts"). Here, we evaluate the few-shot ability of LMs when such held-out examples are unavailable, a setting we call true few-shot learning. We test two model selection criteria, cross-validation and minimum description length, for choosing LM prompts and hyperparameters in the true few-shot setting. On average, both marginally outperform random selection and greatly underperform selection based on held-out examples. Moreover, selection criteria often prefer models that perform significantly worse than randomly-selected ones. We find similar results even when taking into account our uncertainty in a model's true performance during selection, as well as when varying the amount of computation and number of examples used for selection. Overall, our findings suggest that prior work significantly overestimated the true few-shot ability of LMs given the difficulty of few-shot model selection.

[Selective Sampling for Online Best-arm Identification](#)

- Romain Camilleri, Zhihan Xiong, Maryam Fazel, Lalit Jain, Kevin G. Jamieson
- abstract@[open-review](#): This work considers the problem of selective-sampling for best-arm identification. Given a set of potential options $\mathcal{Z} \subset \mathbb{R}^d$, a learner aims to compute with probability greater than $1 - \delta$, $\arg\max_{z \in \mathcal{Z}} z^\top \theta_{\text{last}}$ where θ_{last} is unknown. At each time step, a potential measurement $x_t \in \mathcal{X} \subset \mathbb{R}^d$ is drawn IID and the learner can either choose to take the measurement, in which case they observe a noisy measurement of $x_t^\top \theta_{\text{last}}$, or to abstain from taking the measurement and wait for a potentially more informative point to arrive in the stream. Hence the learner faces a fundamental trade-off between the number of labeled samples they take and when they have collected enough evidence to declare the best arm and stop sampling. The main results of this work precisely characterize this trade-off between labeled samples and stopping time and provide an algorithm that nearly-optimally achieves the minimal label complexity given a desired stopping time. In addition, we show that the optimal decision rule has a simple geometric form based on deciding whether a point is in an ellipse or not. Finally, our framework is general enough to capture binary classification improving upon previous works.

[Multi-task Learning of Order-Consistent Causal Graphs](#)

- Xinshi Chen, Haoran Sun, Caleb Ellington, Eric Xing, Le Song

- abstract@[open-review](#): We consider the problem of discovering K related Gaussian directed acyclic graphs (DAGs), where the involved graph structures share a consistent causal order and sparse unions of supports. Under the multi-task learning setting, we propose a ℓ_1/ℓ_2 -regularized maximum likelihood estimator (MLE) for learning K linear structural equation models. We theoretically show that the joint estimator, by leveraging data across related tasks, can achieve a better sample complexity for recovering the causal order (or topological order) than separate estimations. Moreover, the joint estimator is able to recover non-identifiable DAGs, by estimating them together with some identifiable DAGs. Lastly, our analysis also shows the consistency of union support recovery of the structures. To allow practical implementation, we design a continuous optimization problem whose optimizer is the same as the joint estimator and can be approximated efficiently by an iterative algorithm. We validate the theoretical analysis and the effectiveness of the joint estimator in experiments.

[Learning to Iteratively Solve Routing Problems with Dual-Aspect Collaborative Transformer](#)

- Yining Ma, Jingwen Li, Zhiguang Cao, Wen Song, Le Zhang, Zhenghua Chen, Jing Tang
- abstract@[open-review](#): Recently, Transformer has become a prevailing deep architecture for solving vehicle routing problems (VRPs). However, it is less effective in learning improvement models for VRP because its positional encoding (PE) method is not suitable in representing VRP solutions. This paper presents a novel Dual-Aspect Collaborative Transformer (DACT) to learn embeddings for the node and positional features separately, instead of fusing them together as done in existing ones, so as to avoid potential noises and incompatible correlations. Moreover, the positional features are embedded through a novel cyclic positional encoding (CPE) method to allow Transformer to effectively capture the circularity and symmetry of VRP solutions (i.e., cyclic sequences). We train DACT using Proximal Policy Optimization and design a curriculum learning strategy for better sample efficiency. We apply DACT to solve the traveling salesman problem (TSP) and capacitated vehicle routing problem (CVRP). Results show that our DACT outperforms existing Transformer based improvement models, and exhibits much better generalization performance across different problem sizes on synthetic and benchmark instances, respectively.

[Learning interaction rules from multi-animal trajectories via augmented behavioral models](#)

- Keisuke Fujii, Naoya Takeishi, Kazushi Tsutsui, Emyo Fujioka, Nozomi Nishiumi, Ryoya Tanaka, Mika Fukushiro, Kaoru Ide, Hiroyoshi Kohno, Ken Yoda, Susumu Takahashi, Shizuko Hiryu, Yoshinobu Kawahara
- abstract@[open-review](#): Extracting the interaction rules of biological agents from movement sequences pose challenges in various domains. Granger causality is a practical framework for analyzing the interactions from observed time-series data; however, this framework ignores the structures and assumptions of the generative process in animal behaviors, which may lead to interpretational problems and sometimes erroneous assessments of causality. In this paper, we propose a new framework for learning Granger causality from multi-animal trajectories via augmented theory-based behavioral models with interpretable data-driven models. We adopt an approach for augmenting incomplete multi-agent behavioral models described by time-varying dynamical systems with neural networks. For efficient and interpretable learning, our model leverages theory-based architectures separating navigation and motion processes, and the theory-guided regularization for reliable behavioral modeling. This can provide interpretable signs of Granger-causal effects over time, i.e., when specific others cause the approach or separation. In experiments using synthetic datasets, our method achieved better performance than various baselines. We then analyzed multi-animal datasets of mice, flies, birds, and bats, which verified our method and obtained novel biological insights.

[Differentiable Synthesis of Program Architectures](#)

- Guofeng Cui, He Zhu
- abstract@[open-review](#): Differentiable programs have recently attracted much interest due to their interpretability, compositionality, and their efficiency to leverage differentiable training. However, synthesizing differentiable programs requires optimizing over a combinatorial, rapidly exploded space of program architectures. Despite the development of effective pruning heuristics, previous works essentially enumerate the discrete search space of program architectures, which is inefficient. We propose to encode program architecture search as learning the probability distribution over all possible program derivations induced by a context-free grammar. This allows the search algorithm to efficiently prune away unlikely program derivations to synthesize optimal program architectures. To this end, an efficient gradient-descent based method is developed to conduct program architecture search in a continuous relaxation of the discrete space of grammar rules. Experiment results on four sequence classification tasks demonstrate that our program synthesizer excels in discovering program architectures that lead to differentiable programs with higher F1 scores, while being more efficient than state-of-the-art program synthesis methods.

[Make Sure You're Unsure: A Framework for Verifying Probabilistic Specifications](#)

- Leonard Berrada, Sumanth Dathathri, Krishnamurthy Dvijotham, Robert Stanforth, Rudy R. Bunel, Jonathan Uesato, Sven Gowal, M. Pawan Kumar
- abstract@[open-review](#): Most real world applications require dealing with stochasticity like sensor noise or predictive uncertainty, where formal specifications of desired behavior are inherently probabilistic. Despite the promise of formal verification in ensuring the reliability of neural networks, progress in the direction of probabilistic specifications has been limited. In this direction, we first introduce a general formulation of probabilistic specifications for neural networks, which captures both probabilistic networks (e.g., Bayesian neural networks, MC-Dropout networks) and uncertain inputs (distributions over inputs arising from sensor noise or other perturbations). We then propose a general technique to verify such specifications by generalizing the notion of Lagrangian duality, replacing standard Lagrangian multipliers with "functional multipliers" that can be arbitrary functions of the activations at a given layer. We show that an optimal choice of functional multipliers leads to exact verification (i.e., sound and complete verification), and for specific forms of multipliers, we develop tractable practical verification algorithms. We empirically validate our algorithms by applying them to Bayesian Neural Networks (BNNs) and MC Dropout Networks, and certifying properties such as adversarial robustness and robust detection of out-of-distribution (OOD) data. On these tasks we are able to provide significantly stronger guarantees when compared to prior work -- for instance, for a VGG-64 MC-Dropout CNN trained on CIFAR-10 in a verification-agnostic manner, we improve the certified AUC (a verified lower bound on the true AUC) for robust OOD detection (on CIFAR-100) from 0% to 29% . Similarly, for a BNN trained on MNIST, we improve on the ℓ_∞ robust accuracy from 60.2% to 74.6% . Further, on a novel specification -- distributionally robust OOD detection -- we improve on the certified AUC from 5% to 23% .

[Oracle-Efficient Regret Minimization in Factored MDPs with Unknown Structure](#)

- Aviv Rosenberg, Yishay Mansour
- abstract@[open-review](#): We study regret minimization in non-episodic factored Markov decision processes (FMDPs), where all existing algorithms make the strong assumption that the factored structure of the FMDP is known to the learner in advance. In this paper, we provide the first algorithm that learns the structure of the FMDP while minimizing the regret. Our algorithm is based on the optimism in face of uncertainty principle, combined with a simple statistical method for structure learning, and can be implemented efficiently given oracle-access to an FMDP planner. Moreover, we give a variant of our algorithm that remains efficient even when the oracle is limited to non-factored actions, which is the case with almost all existing approximate planners. Finally, we leverage our techniques to prove a novel lower bound for the known structure case, closing the gap to the regret bound of Chen et al. [2021].

[Linear-Time Probabilistic Solution of Boundary Value Problems](#)

- Nicholas Krämer, Philipp Hennig

- abstract@[open-review](#): We propose a fast algorithm for the probabilistic solution of boundary value problems (BVPs), which are ordinary differential equations subject to boundary conditions. In contrast to previous work, we introduce a Gauss-Markov prior and tailor it specifically to BVPs, which allows computing a posterior distribution over the solution in linear time, at a quality and cost comparable to that of well-established, non-probabilistic methods. Our model further delivers uncertainty quantification, mesh refinement, and hyperparameter adaptation. We demonstrate how these practical considerations positively impact the efficiency of the scheme. Altogether, this results in a practically usable probabilistic BVP solver that is (in contrast to non-probabilistic algorithms) natively compatible with other parts of the statistical modelling tool-chain.

[Lifelong Domain Adaptation via Consolidated Internal Distribution](#)

- Mohammad Rostami
- abstract@[open-review](#): We develop an algorithm to address unsupervised domain adaptation (UDA) in continual learning (CL) settings. The goal is to update a model continually to learn distributional shifts across sequentially arriving tasks with unlabeled data while retaining the knowledge about the past learned tasks. Existing UDA algorithms address the challenge of domain shift, but they require simultaneous access to the datasets of the source and the target domains. On the other hand, existing works on CL can handle tasks with labeled data. Our solution is based on consolidating the learned internal distribution for improved model generalization on new domains and benefitting from experience replay to overcome catastrophic forgetting.

[Counterbalancing Learning and Strategic Incentives in Allocation Markets](#)

- Jamie Kang, Faidra Monachou, Moran Koren, Itai Ashlagi
- abstract@[open-review](#): Motivated by the high discard rate of donated organs in the United States, we study an allocation problem in the presence of learning and strategic incentives. We consider a setting where a benevolent social planner decides whether and how to allocate a single indivisible object to a queue of strategic agents. The object has a common true quality, good or bad, which is ex-ante unknown to everyone. Each agent holds an informative, yet noisy, private signal about the quality. To make a correct allocation decision the planner attempts to learn the object quality by truthfully eliciting agents' signals. Under the commonly applied sequential offering mechanism, we show that learning is hampered by the presence of strategic incentives as herding may emerge. This can result in incorrect allocation and welfare loss. To overcome these issues, we propose a novel class of incentive-compatible mechanisms. Our mechanism involves a batch-by-batch, dynamic voting process using a majority rule. We prove that the proposed voting mechanisms improve the probability of correct allocation whenever agents are sufficiently well informed. Particularly, we show that such an improvement can be achieved via a simple greedy algorithm. We quantify the improvement using simulations.

[Controlling Neural Networks with Rule Representations](#)

- Sungyong Seo, Sercan Arik, Jinsung Yoon, Xiang Zhang, Kihyuk Sohn, Tomas Pfister
- abstract@[open-review](#): We propose a novel training method that integrates rules into deep learning, in a way the strengths of the rules are controllable at inference. Deep Neural Networks with Controllable Rule Representations (DeepCTRL) incorporates a rule encoder into the model coupled with a rule-based objective, enabling a shared representation for decision making. DeepCTRL is agnostic to data type and model architecture. It can be applied to any kind of rule defined for inputs and outputs. The key aspect of DeepCTRL is that it does not require retraining to adapt the rule strength -- at inference, the user can adjust it based on the desired operation point on accuracy vs. rule verification ratio. In real-world domains where incorporating rules is critical -- such as Physics, Retail and Healthcare -- we show the effectiveness of DeepCTRL in teaching rules for deep learning. DeepCTRL improves the trust and reliability of the trained models by significantly increasing their rule verification ratio, while also providing accuracy gains at downstream tasks. Additionally, DeepCTRL enables novel use cases such as hypothesis testing of the rules on data samples, and unsupervised adaptation based on shared rules between datasets.

[Making the most of your day: online learning for optimal allocation of time](#)

- Etienne Boursier, Tristan Garrec, Vianney Perchet, Marco Scarsini
- abstract@[open-review](#): We study online learning for optimal allocation when the resource to be allocated is time. An agent receives task proposals sequentially according to a Poisson process and can either accept or reject a proposed task. If she accepts the proposal, she is busy for the duration of the task and obtains a reward that depends on the task duration. If she rejects it, she remains on hold until a new task proposal arrives. We study the regret incurred by the agent first when she knows her reward function but does not know the distribution of the task duration, and then when she does not know her reward function, either. Faster rates are finally obtained by adding structural assumptions on the distribution of rides or on the reward function. This natural setting bears similarities with contextual (one-armed) bandits, but with the crucial difference that the normalized reward associated to a context depends on the whole distribution of contexts.

[Federated Reconstruction: Partially Local Federated Learning](#)

- Karan Singh, Hakim Sidahmed, Zachary Garrett, Shanshan Wu, John Rush, Sushant Prakash
- abstract@[open-review](#): Personalization methods in federated learning aim to balance the benefits of federated and local training for data availability, communication cost, and robustness to client heterogeneity. Approaches that require clients to communicate all model parameters can be undesirable due to privacy and communication constraints. Other approaches require always-available or stateful clients, impractical in large-scale cross-device settings. We introduce Federated Reconstruction, the first model-agnostic framework for partially local federated learning suitable for training and inference at scale. We motivate the framework via a connection to model-agnostic meta learning, empirically demonstrate its performance over existing approaches for collaborative filtering and next word prediction, and release an open-source library for evaluating approaches in this setting. We also describe the successful deployment of this approach at scale for federated collaborative filtering in a mobile keyboard application.

[Optimal prediction of Markov chains with and without spectral gap](#)

- Yanjun Han, Soham Jana, Yihong Wu
- abstract@[open-review](#): We study the following learning problem with dependent data: Given a trajectory of length n from a stationary Markov chain with k states, the goal is to predict the distribution of the next state. For $k \leq O(\sqrt{n})$, the optimal prediction risk in the Kullback-Leibler divergence is shown to be $\Theta(\frac{1}{n} \log \frac{n}{k^2})$, in contrast to the optimal rate of $\Theta(\log \log n)$ for $k=2$ previously shown in Falahatgar et al in 2016. These nonparametric rates can be attributed to the memory in the data, as the spectral gap of the Markov chain can be arbitrarily small. To quantify the memory effect, we study irreducible reversible chains with a prescribed spectral gap. In addition to characterizing the optimal prediction risk for two states, we show that, as long as the spectral gap is not excessively small, the prediction risk in the Markov model is $O(\frac{1}{n})$, which coincides with that of an iid model with the same number of parameters.

[Subquadratic Overparameterization for Shallow Neural Networks](#)

- ChaeHwan Song, Ali Ramezani-Kebrya, Thomas Pethick, Armin Eftekhari, Volkan Cevher
- abstract@[open-review](#): Overparameterization refers to the important phenomenon where the width of a neural network is chosen such that learning algorithms can provably attain zero loss in nonconvex training. The existing theory establishes such global convergence using various initialization strategies, training modifications, and width scalings. In particular, the state-of-the-art results require the width to scale quadratically with the number of

training data under standard initialization strategies used in practice for best generalization performance. In contrast, the most recent results obtain linear scaling either with requiring initializations that lead to the "lazy-training", or training only a single layer. In this work, we provide an analytical framework that allows us to adopt standard initialization strategies, possibly avoid lazy training, and train all layers simultaneously in basic shallow neural networks while attaining a desirable subquadratic scaling on the network width. We achieve the desiderata via Polyak-Lojasiewicz condition, smoothness, and standard assumptions on data, and use tools from random matrix theory.

[Continuous Doubly Constrained Batch Reinforcement Learning](#)

- Rasool Fakoor, Jonas W. Mueller, Kavosh Asadi, Pratik Chaudhari, Alexander J. Smola
- abstract@[open-review](#): Reliant on too many experiments to learn good actions, current Reinforcement Learning (RL) algorithms have limited applicability in real-world settings, which can be too expensive to allow exploration. We propose an algorithm for batch RL, where effective policies are learned using only a fixed offline dataset instead of online interactions with the environment. The limited data in batch RL produces inherent uncertainty in value estimates of states/actions that were insufficiently represented in the training data. This leads to particularly severe extrapolation when our candidate policies diverge from one that generated the data. We propose to mitigate this issue via two straightforward penalties: a policy-constraint to reduce this divergence and a value-constraint that discourages overly optimistic estimates. Over a comprehensive set of \$32\\$ continuous-action batch RL benchmarks, our approach compares favorably to state-of-the-art methods, regardless of how the offline data were collected.

[Bridging Explicit and Implicit Deep Generative Models via Neural Stein Estimators](#)

- Qitian Wu, Rui Gao, Hongyuan Zha
- abstract@[open-review](#): There are two types of deep generative models: explicit and implicit. The former defines an explicit density form that allows likelihood inference; while the latter targets a flexible transformation from random noise to generated samples. While the two classes of generative models have shown great power in many applications, both of them, when used alone, suffer from respective limitations and drawbacks. To take full advantages of both models and enable mutual compensation, we propose a novel joint training framework that bridges an explicit (unnormalized) density estimator and an implicit sample generator via Stein discrepancy. We show that our method 1) induces novel mutual regularization via kernel Sobolev norm penalization and Moreau-Yosida regularization, and 2) stabilizes the training dynamics. Empirically, we demonstrate that proposed method can facilitate the density estimator to more accurately identify data modes and guide the generator to output higher-quality samples, comparing with training a single counterpart. The new approach also shows promising results when the training samples are contaminated or limited.

[Score-based Generative Modeling in Latent Space](#)

- Arash Vahdat, Karsten Kreis, Jan Kautz
- abstract@[open-review](#): Score-based generative models (SGMs) have recently demonstrated impressive results in terms of both sample quality and distribution coverage. However, they are usually applied directly in data space and often require thousands of network evaluations for sampling. Here, we propose the Latent Score-based Generative Model (LSGM), a novel approach that trains SGMs in a latent space, relying on the variational autoencoder framework. Moving from data to latent space allows us to train more expressive generative models, apply SGMs to non-continuous data, and learn smoother SGMs in a smaller space, resulting in fewer network evaluations and faster sampling. To enable training LSGMs end-to-end in a scalable and stable manner, we (i) introduce a new score-matching objective suitable to the LSGM setting, (ii) propose a novel parameterization of the score function that allows SGM to focus on the mismatch of the target distribution with respect to a simple Normal one, and (iii) analytically derive multiple techniques for variance reduction of the training objective. LSGM obtains a state-of-the-art FID score of 2.10 on CIFAR-10, outperforming all existing generative results on this dataset. On CelebA-HQ-256, LSGM is on a par with previous SGMs in sample quality while outperforming them in sampling time by two orders of magnitude. In modeling binary images, LSGM achieves state-of-the-art likelihood on the binarized OMNIGLOT dataset.

[Deep Conditional Gaussian Mixture Model for Constrained Clustering](#)

- Laura Manduchi, Kieran Chin-Cheong, Holger Michel, Sven Wellmann, Julia Vogt
- abstract@[open-review](#): Constrained clustering has gained significant attention in the field of machine learning as it can leverage prior information on a growing amount of only partially labeled data. Following recent advances in deep generative models, we propose a novel framework for constrained clustering that is intuitive, interpretable, and can be trained efficiently in the framework of stochastic gradient variational inference. By explicitly integrating domain knowledge in the form of probabilistic relations, our proposed model (DC-GMM) uncovers the underlying distribution of data conditioned on prior clustering preferences, expressed as \textit{pairwise constraints}. These constraints guide the clustering process towards a desirable partition of the data by indicating which samples should or should not belong to the same cluster. We provide extensive experiments to demonstrate that DC-GMM shows superior clustering performances and robustness compared to state-of-the-art deep constrained clustering methods on a wide range of data sets. We further demonstrate the usefulness of our approach on two challenging real-world applications.

[Bootstrap Your Object Detector via Mixed Training](#)

- Mengde Xu, Zheng Zhang, Fangyun Wei, Yutong Lin, Yue Cao, Stephen Lin, Han Hu, Xiang Bai
- abstract@[open-review](#): We introduce MixTraining, a new training paradigm for object detection that can improve the performance of existing detectors for free. MixTraining enhances data augmentation by utilizing augmentations of different strengths while excluding the strong augmentations of certain training samples that may be detrimental to training. In addition, it addresses localization noise and missing labels in human annotations by incorporating pseudo boxes that can compensate for these errors. Both of these MixTraining capabilities are made possible through bootstrapping on the detector, which can be used to predict the difficulty of training on a strong augmentation, as well as to generate reliable pseudo boxes thanks to the robustness of neural networks to labeling error. MixTraining is found to bring consistent improvements across various detectors on the COCO dataset. In particular, the performance of Faster R-CNN~\cite{ren2015faster} with a ResNet-50~\cite{he2016deep} backbone is improved from 41.7 mAP to 44.0 mAP, and the accuracy of Cascade-RCNN~\cite{cai2018cascade} with a Swin-Small~\cite{liu2021swin} backbone is raised from 50.9 mAP to 52.8 mAP.

[Tensor decompositions of higher-order correlations by nonlinear Hebbian plasticity](#)

- Gabriel Ocker, Michael Buice
- abstract@[open-review](#): Biological synaptic plasticity exhibits nonlinearities that are not accounted for by classic Hebbian learning rules. Here, we introduce a simple family of generalized nonlinear Hebbian learning rules. We study the computations implemented by their dynamics in the simple setting of a neuron receiving feedforward inputs. These nonlinear Hebbian rules allow a neuron to learn tensor decompositions of its higher-order input correlations. The particular input correlation decomposed and the form of the decomposition depend on the location of nonlinearities in the plasticity rule. For simple, biologically motivated parameters, the neuron learns eigenvectors of higher-order input correlation tensors. We prove that tensor eigenvectors are attractors and determine their basins of attraction. We calculate the volume of those basins, showing that the dominant eigenvector has the largest basin of attraction. We then study arbitrary learning rules and find that any learning rule that admits a finite Taylor expansion into the neural input and output also has stable equilibria at generalized eigenvectors of higher-order input correlation tensors. Nonlinearities in synaptic plasticity thus allow a neuron to encode higher-order input correlations in a simple fashion.

[Online Adaptation to Label Distribution Shift](#)

- Ruihan Wu, Chuan Guo, Yi Su, Kilian Q. Weinberger
- abstract@[open-review](#): Machine learning models often encounter distribution shifts when deployed in the real world. In this paper, we focus on adaptation to label distribution shift in the online setting, where the test-time label distribution is continually changing and the model must dynamically adapt to it without observing the true label. This setting is common in many real world scenarios such as medical diagnosis, where disease prevalences can vary substantially at different times of the year. Leveraging a novel analysis, we show that the lack of true label does not hinder estimation of the expected test loss, which enables the reduction of online label shift adaptation to conventional online learning. Informed by this observation, we propose adaptation algorithms inspired by classical online learning techniques such as Follow The Leader (FTL) and Online Gradient Descent (OGD) and derive their regret bounds. We empirically verify our findings under both simulated and real world label distribution shifts and show that OGD is particularly effective and robust to a variety of challenging label shift scenarios.

[**One Explanation is Not Enough: Structured Attention Graphs for Image Classification**](#)

- Vivswan Shitole, Fuxin Li, Minsuk Kahng, Prasad Tadepalli, Alan Fern
- abstract@[open-review](#): Attention maps are popular tools for explaining the decisions of convolutional neural networks (CNNs) for image classification. Typically, for each image of interest, a single attention map is produced, which assigns weights to pixels based on their importance to the classification. We argue that a single attention map provides an incomplete understanding since there are often many other maps that explain a classification equally well. In this paper, we propose to utilize a beam search algorithm to systematically search for multiple explanations for each image. Results show that there are indeed multiple relatively localized explanations for many images. However, naively showing multiple explanations to users can be overwhelming and does not reveal their common and distinct structures. We introduce structured attention graphs (SAGs), which compactly represent sets of attention maps for an image by visualizing how different combinations of image regions impact the confidence of a classifier. An approach to computing a compact and representative SAG for visualization is proposed via diverse sampling. We conduct a user study comparing the use of SAGs to traditional attention maps for answering comparative counterfactual questions about image classifications. Our results show that the users are significantly more accurate when presented with SAGs compared to standard attention map baselines.

[**Integrating Expert ODEs into Neural ODEs: Pharmacology and Disease Progression**](#)

- Zhaozhi Qian, William Zame, Lucas Fleuren, Paul Elbers, Mihaela van der Schaar
- abstract@[open-review](#): Modeling a system's temporal behaviour in reaction to external stimuli is a fundamental problem in many areas. Pure Machine Learning (ML) approaches often fail in the small sample regime and cannot provide actionable insights beyond predictions. A promising modification has been to incorporate expert domain knowledge into ML models. The application we consider is predicting the patient health status and disease progression over time, where a wealth of domain knowledge is available from pharmacology. Pharmacological models describe the dynamics of carefully-chosen medically meaningful variables in terms of systems of Ordinary Differential Equations (ODEs). However, these models only describe a limited collection of variables, and these variables are often not observable in clinical environments. To close this gap, we propose the latent hybridisation model (LHM) that integrates a system of expert-designed ODEs with machine-learned Neural ODEs to fully describe the dynamics of the system and to link the expert and latent variables to observable quantities. We evaluated LHM on synthetic data as well as real-world intensive care data of COVID-19 patients. LHM consistently outperforms previous works, especially when few training samples are available such as at the beginning of the pandemic.

[**Shifted Chunk Transformer for Spatio-Temporal Representational Learning**](#)

- Xuefan Zha, Wentao Zhu, Lv Xun, Sen Yang, Ji Liu
- abstract@[open-review](#): Spatio-temporal representational learning has been widely adopted in various fields such as action recognition, video object segmentation, and action anticipation. Previous spatio-temporal representational learning approaches primarily employ ConvNets or sequential models, e.g., LSTM, to learn the intra-frame and inter-frame features. Recently, Transformer models have successfully dominated the study of natural language processing (NLP), image classification, etc. However, the pure-Transformer based spatio-temporal learning can be prohibitively costly on memory and computation to extract fine-grained features from a tiny patch. To tackle the training difficulty and enhance the spatio-temporal learning, we construct a shifted chunk Transformer with pure self-attention blocks. Leveraging the recent efficient Transformer design in NLP, this shifted chunk Transformer can learn hierarchical spatio-temporal features from a local tiny patch to a global videoclip. Our shifted self-attention can also effectively model complicated inter-frame variances. Furthermore, we build a clip encoder based on Transformer to model long-term temporal dependencies. We conduct thorough ablation studies to validate each component and hyper-parameters in our shifted chunk Transformer, and it outperforms previous state-of-the-art approaches on Kinetics-400, Kinetics-600, UCF101, and HMDB51.

[**Faster proximal algorithms for matrix optimization using Jacobi-based eigenvalue methods**](#)

- Hamza Fawzi, Harry Goulbourne
- abstract@[open-review](#): We consider proximal splitting algorithms for convex optimization problems over matrices. A significant computational bottleneck in many of these algorithms is the need to compute a full eigenvalue or singular value decomposition at each iteration for the evaluation of a proximal operator. In this paper we propose to use an old and surprisingly simple method due to Jacobi to compute these eigenvalue and singular value decompositions, and we demonstrate that it can lead to substantial gains in terms of computation time compared to standard approaches. We rely on three essential properties of this method: (a) its ability to exploit an approximate decomposition as an initial point, which in the case of iterative optimization algorithms can be obtained from the previous iterate; (b) its parallel nature which makes it a great fit for hardware accelerators such as GPUs, now common in machine learning, and (c) its simple termination criterion which allows us to trade-off accuracy with computation time. We demonstrate the efficacy of this approach on a variety of algorithms and problems, and show that, on a GPU, we can obtain 5 to 10x speed-ups in the evaluation of proximal operators compared to standard CPU or GPU linear algebra routines. Our findings are supported by new theoretical results providing guarantees on the approximation quality of proximal operators obtained using approximate eigenvalue or singular value decompositions.

[**Decrypting Cryptic Crosswords: Semantically Complex Wordplay Puzzles as a Target for NLP**](#)

- Josh Rozner, Christopher Potts, Kyle Mahowald
- abstract@[open-review](#): Cryptic crosswords, the dominant crossword variety in the UK, are a promising target for advancing NLP systems that seek to process semantically complex, highly compositional language. Cryptic clues read like fluent natural language but are adversarially composed of two parts: a definition and a wordplay cipher requiring character-level manipulations. Expert humans use creative intelligence to solve cryptics, flexibly combining linguistic, world, and domain knowledge. In this paper, we make two main contributions. First, we present a dataset of cryptic clues as a challenging new benchmark for NLP systems that seek to process compositional language in more creative, human-like ways. After showing that three non-neural approaches and T5, a state-of-the-art neural language model, do not achieve good performance, we make our second main contribution: a novel curriculum approach, in which the model is first fine-tuned on related tasks such as unscrambling words. We also introduce a challenging data split, examine the meta-linguistic capabilities of subword-tokenized models, and investigate model systematicity by perturbing the wordplay part of clues, showing that T5 exhibits behavior partially consistent with human solving strategies. Although our curricular approach considerably improves on the T5 baseline, our best-performing model still fails to generalize to the extent that humans can. Thus, cryptic crosswords remain an unsolved challenge for NLP systems and a potential source of future innovation.

[**An Improved Analysis of Gradient Tracking for Decentralized Machine Learning**](#)

- Anastasiia Koloskova, Tao Lin, Sebastian U. Stich
- abstract@[open-review](#): We consider decentralized machine learning over a network where the training data is distributed across n agents, each of which can compute stochastic model updates on their local data. The agent's common goal is to find a model that minimizes the average of all local loss functions. While gradient tracking (GT) algorithms can overcome a key challenge, namely accounting for differences between workers' local data distributions, the known convergence rates for GT algorithms are not optimal with respect to their dependence on the mixing parameter p (related to the spectral gap of the connectivity matrix). We provide a tighter analysis of the GT method in the stochastic strongly convex, convex and non-convex settings. We improve the dependency on p from $\mathcal{O}(p^{-2})$ to $\mathcal{O}(p^{-1}c^{-1})$ in the noiseless case and from $\mathcal{O}(p^{-3/2})$ to $\mathcal{O}(p^{-1/2}c^{-1})$ in the general stochastic case, where $c \geq p$ is related to the negative eigenvalues of the connectivity matrix (and is a constant in most practical applications). This improvement was possible due to a new proof technique which could be of independent interest.

[Entropic Desired Dynamics for Intrinsic Control](#)

- Steven Hansen, Guillaume Desjardins, Kate Baumli, David Warde-Farley, Nicolas Heess, Simon Osindero, Volodymyr Mnih
- abstract@[open-review](#): An agent might be said, informally, to have mastery of its environment when it has maximised the effective number of states it can reliably reach. In practice, this often means maximizing the number of latent codes that can be discriminated from future states under some short time horizon (e.g. \cite{eysenbach2018diversity}). By situating these latent codes in a globally consistent coordinate system, we show that agents can reliably reach more states in the long term while still optimizing a local objective. A simple instantiation of this idea, Entropic Desired Dynamics for Intrinsic Control (EDDICT), assumes fixed additive latent dynamics, which results in tractable learning and an interpretable latent space. Compared to prior methods, EDDICT's globally consistent codes allow it to be far more exploratory, as demonstrated by improved state coverage and increased unsupervised performance on hard exploration games such as Montezuma's Revenge.

[Exploring Cross-Video and Cross-Modality Signals for Weakly-Supervised Audio-Visual Video Parsing](#)

- Yan-Bo Lin, Hung-Yu Tseng, Hsin-Ying Lee, Yen-Yu Lin, Ming-Hsuan Yang
- abstract@[open-review](#): The audio-visual video parsing task aims to temporally parse a video into audio or visual event categories. However, it is labor intensive to temporally annotate audio and visual events and thus hampers the learning of a parsing model. To this end, we propose to explore additional cross-video and cross-modality supervisory signals to facilitate weakly-supervised audio-visual video parsing. The proposed method exploits both the common and diverse event semantics across videos to identify audio or visual events. In addition, our method explores event co-occurrence across audio, visual, and audio-visual streams. We leverage the explored cross-modality co-occurrence to localize segments of target events while excluding irrelevant ones. The discovered supervisory signals across different videos and modalities can greatly facilitate the training with only video-level annotations. Quantitative and qualitative results demonstrate that the proposed method performs favorably against existing methods on weakly-supervised audio-visual video parsing.

[Littlestone Classes are Privately Online Learnable](#)

- Noah Golowich, Roi Livni
- abstract@[open-review](#): We consider the problem of online classification under a privacy constraint. In this setting a learner observes sequentially a stream of labelled examples (x_t, y_t) , for $1 \leq t \leq T$, and returns at each iteration t a hypothesis h_t which is used to predict the label of each new example x_t . The learner's performance is measured by her regret against a known hypothesis class \mathcal{H} . We require that the algorithm satisfies the following privacy constraint: the sequence h_1, \dots, h_T of hypotheses output by the algorithm needs to be an (ϵ, δ) -differentially private function of the whole input sequence $(x_1, y_1), \dots, (x_T, y_T)$. We provide the first non-trivial regret bound for the realizable setting. Specifically, we show that if the class \mathcal{H} has constant Littlestone dimension then, given an oblivious sequence of labelled examples, there is a private learner that makes in expectation at most $O(\log T)$ mistakes -- comparable to the optimal mistake bound in the non-private case, up to a logarithmic factor. Moreover, for general values of the Littlestone dimension d , the same mistake bound holds but with a doubly-exponential in d factor. A recent line of work has demonstrated a strong connection between classes that are online learnable and those that are differentially-private learnable. Our results strengthen this connection and show that an online learning algorithm can in fact be directly privatized (in the realizable setting). We also discuss an adaptive setting and provide a sublinear regret bound of $O(\sqrt{T})$.

[Dual Parameterization of Sparse Variational Gaussian Processes](#)

- Vincent ADAM, Paul Chang, Mohammad Emtiyaz E. Khan, Arno Solin
- abstract@[open-review](#): Sparse variational Gaussian process (SVGP) methods are a common choice for non-conjugate Gaussian process inference because of their computational benefits. In this paper, we improve their computational efficiency by using a dual parameterization where each data example is assigned dual parameters, similarly to site parameters used in expectation propagation. Our dual parameterization speeds-up inference using natural gradient descent, and provides a tighter evidence lower bound for hyperparameter learning. The approach has the same memory cost as the current SVGP methods, but it is faster and more accurate.

[Learning to dehaze with polarization](#)

- Chu Zhou, Minggui Teng, Yufei Han, Chao Xu, Boxin Shi
- abstract@[open-review](#): Haze, a common kind of bad weather caused by atmospheric scattering, decreases the visibility of scenes and degenerates the performance of computer vision algorithms. Single-image dehazing methods have shown their effectiveness in a large variety of scenes, however, they are based on handcrafted priors or learned features, which do not generalize well to real-world images. Polarization information can be used to relieve its ill-posedness, however, real-world images are still challenging since existing polarization-based methods usually assume that the transmitted light is not significantly polarized, and they require specific clues to estimate necessary physical parameters. In this paper, we propose a generalized physical formation model of hazy images and a robust polarization-based dehazing pipeline without the above assumption or requirement, along with a neural network tailored to the pipeline. Experimental results show that our approach achieves state-of-the-art performance on both synthetic data and real-world hazy images.

[Conservative Data Sharing for Multi-Task Offline Reinforcement Learning](#)

- Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Sergey Levine, Chelsea Finn
- abstract@[open-review](#): Offline reinforcement learning (RL) algorithms have shown promising results in domains where abundant pre-collected data is available. However, prior methods focus on solving individual problems from scratch with an offline dataset without considering how an offline RL agent can acquire multiple skills. We argue that a natural use case of offline RL is in settings where we can pool large amounts of data collected in various scenarios for solving different tasks, and utilize all of this data to learn behaviors for all the tasks more effectively rather than training each one in isolation. However, sharing data across all tasks in multi-task offline RL performs surprisingly poorly in practice. Thorough empirical analysis, we find that sharing data can actually exacerbate the distributional shift between the learned policy and the dataset, which in turn can lead to divergence of the learned policy and poor performance. To address this challenge, we develop a simple technique for data-sharing in multi-task offline RL that routes data based on the improvement over the task-specific data. We call this approach conservative data sharing (CDS), and it can be applied with multiple single-

task offline RL methods. On a range of challenging multi-task locomotion, navigation, and vision-based robotic manipulation problems, CDS achieves the best or comparable performance compared to prior offline multi-task RL methods and previous data sharing approaches.

[Universal Rate-Distortion-Perception Representations for Lossy Compression](#)

- George Zhang, Jingjing Qian, Jun Chen, Ashish Khisti
- abstract@[open-review](#): In the context of lossy compression, Blau \& Michaeli (2019) adopt a mathematical notion of perceptual quality and define the information rate-distortion-perception function, generalizing the classical rate-distortion tradeoff. We consider the notion of universal representations in which one may fix an encoder and vary the decoder to achieve any point within a collection of distortion and perception constraints. We prove that the corresponding information-theoretic universal rate-distortion-perception function is operationally achievable in an approximate sense. Under MSE distortion, we show that the entire distortion-perception tradeoff of a Gaussian source can be achieved by a single encoder of the same rate asymptotically. We then characterize the achievable distortion-perception region for a fixed representation in the case of arbitrary distributions, and identify conditions under which the aforementioned results continue to hold approximately. This motivates the study of practical constructions that are approximately universal across the RDP tradeoff, thereby alleviating the need to design a new encoder for each objective. We provide experimental results on MNIST and SVHN suggesting that on image compression tasks, the operational tradeoffs achieved by machine learning models with a fixed encoder suffer only a small penalty when compared to their variable encoder counterparts.

[What's a good imputation to predict with missing values?](#)

- Marine Le Morvan, Julie Josse, Erwan Scornet, Gael Varoquaux
- abstract@[open-review](#): How to learn a good predictor on data with missing values? Most efforts focus on first imputing as well as possible and second learning on the completed data to predict the outcome. Yet, this widespread practice has no theoretical grounding. Here we show that for almost all imputation functions, an impute-then-regress procedure with a powerful learner is Bayes optimal. This result holds for all missing-values mechanisms, in contrast with the classic statistical results that require missing-at-random settings to use imputation in probabilistic modeling. Moreover, it implies that perfect conditional imputation is not needed for good prediction asymptotically. In fact, we show that on perfectly imputed data the best regression function will generally be discontinuous, which makes it hard to learn. Crafting instead the imputation so as to leave the regression function unchanged simply shifts the problem to learning discontinuous imputations. Rather, we suggest that it is easier to learn imputation and regression jointly. We propose such a procedure, adapting NeuMiss, a neural network capturing the conditional links across observed and unobserved variables whatever the missing-value pattern. Our experiments confirm that joint imputation and regression through NeuMiss is better than various two step procedures in a finite-sample regime.

[Replacing Rewards with Examples: Example-Based Policy Search via Recursive Classification](#)

- Ben Eysenbach, Sergey Levine, Russ R. Salakhutdinov
- abstract@[open-review](#): Reinforcement learning (RL) algorithms assume that users specify tasks by manually writing down a reward function. However, this process can be laborious and demands considerable technical expertise. Can we devise RL algorithms that instead enable users to specify tasks simply by providing examples of successful outcomes? In this paper, we derive a control algorithm that maximizes the future probability of these successful outcome examples. Prior work has approached similar problems with a two-stage process, first learning a reward function and then optimizing this reward function using another reinforcement learning algorithm. In contrast, our method directly learns a value function from transitions and successful outcomes, without learning this intermediate reward function. Our method therefore requires fewer hyperparameters to tune and lines of code to debug. We show that our method satisfies a new data-driven Bellman equation, where examples take the place of the typical reward function term. Experiments show that our approach outperforms prior methods that learn explicit reward functions.

[Hierarchical Skills for Efficient Exploration](#)

- Jonas Gehring, Gabriel Synnaeve, Andreas Krause, Nicolas Usunier
- abstract@[open-review](#): In reinforcement learning, pre-trained low-level skills have the potential to greatly facilitate exploration. However, prior knowledge of the downstream task is required to strike the right balance between generality (fine-grained control) and specificity (faster learning) in skill design. In previous work on continuous control, the sensitivity of methods to this trade-off has not been addressed explicitly, as locomotion provides a suitable prior for navigation tasks, which have been of foremost interest. In this work, we analyze this trade-off for low-level policy pre-training with a new benchmark suite of diverse, sparse-reward tasks for bipedal robots. We alleviate the need for prior knowledge by proposing a hierarchical skill learning framework that acquires skills of varying complexity in an unsupervised manner. For utilization on downstream tasks, we present a three-layered hierarchical learning algorithm to automatically trade off between general and specific skills as required by the respective task. In our experiments, we show that our approach performs this trade-off effectively and achieves better results than current state-of-the-art methods for end-to-end hierarchical reinforcement learning and unsupervised skill discovery.

[Evidential Softmax for Sparse Multimodal Distributions in Deep Generative Models](#)

- Phil Chen, Mikhaylo Itkina, Ransalu Senanayake, Mykel J Kochenderfer
- abstract@[open-review](#): Many applications of generative models rely on the marginalization of their high-dimensional output probability distributions. Normalization functions that yield sparse probability distributions can make exact marginalization more computationally tractable. However, sparse normalization functions usually require alternative loss functions for training since the log-likelihood is undefined for sparse probability distributions. Furthermore, many sparse normalization functions often collapse the multimodality of distributions. In this work, we present ev-softmax, a sparse normalization function that preserves the multimodality of probability distributions. We derive its properties, including its gradient in closed-form, and introduce a continuous family of approximations to ev-softmax that have full support and can be trained with probabilistic loss functions such as negative log-likelihood and Kullback-Leibler divergence. We evaluate our method on a variety of generative models, including variational autoencoders and autoregressive architectures. Our method outperforms existing dense and sparse normalization techniques in distributional accuracy. We demonstrate that ev-softmax successfully reduces the dimensionality of probability distributions while maintaining multimodality.

[Submodular + Concave](#)

- Siddharth Mitra, Moran Feldman, Amin Karbasi
- abstract@[open-review](#): It has been well established that first order optimization methods can converge to the maximal objective value of concave functions and provide constant factor approximation guarantees for (non-convex/non-concave) continuous submodular functions. In this work, we initiate the study of the maximization of functions of the form $\$F(x) = G(x) + C(x)\$$ over a solvable convex body $\$P\$$, where $\$G\$$ is a smooth DR-submodular function and $\$C\$$ is a smooth concave function. This class of functions is a strict extension of both concave and continuous DR-submodular functions for which no theoretical guarantee is known. We provide a suite of Frank-Wolfe style algorithms, which, depending on the nature of the objective function (i.e., if $\$G\$$ and $\$C\$$ are monotone or not, and non-negative or not) and on the nature of the set $\$P\$$ (i.e., whether it is downward closed or not), provide $\$1/1/e\$, \$1/e\$,$ or $\$1/2\$$ approximation guarantees. We then use our algorithms to get a framework to smoothly interpolate between choosing a diverse set of elements from a given ground set (corresponding to the mode of a determinantal point process) and choosing a clustered set of elements (corresponding to the maxima of a suitable concave function). Additionally, we apply our algorithms to various functions in the above class (DR-submodular + concave) in both constrained and unconstrained settings, and show that our algorithms consistently outperform natural baselines.

[DeepGEM: Generalized Expectation-Maximization for Blind Inversion](#)

- Angela Gao, Jorge Castellanos, Yisong Yue, Zachary Ross, Katherine Bouman
- abstract@[open-review](#): Typically, inversion algorithms assume that a forward model, which relates a source to its resulting measurements, is known and fixed. Using collected indirect measurements and the forward model, the goal becomes to recover the source. When the forward model is unknown, or imperfect, artifacts due to model mismatch occur in the recovery of the source. In this paper, we study the problem of blind inversion: solving an inverse problem with unknown or imperfect knowledge of the forward model parameters. We propose DeepGEM, a variational Expectation-Maximization (EM) framework that can be used to solve for the unknown parameters of the forward model in an unsupervised manner. DeepGEM makes use of a normalizing flow generative network to efficiently capture complex posterior distributions, which leads to more accurate evaluation of the source's posterior distribution used in EM. We showcase the effectiveness of our DeepGEM approach by achieving strong performance on the challenging problem of blind seismic tomography, where we significantly outperform the standard method used in seismology. We also demonstrate the generality of DeepGEM by applying it to a simple case of blind deconvolution.

[Learning to Generate Visual Questions with Noisy Supervision](#)

- Shen Kai, Lingfei Wu, Siliang Tang, Yueling Zhuang, zhen he, Zhuoye Ding, Yun Xiao, Bo Long
- abstract@[open-review](#): The task of visual question generation (VQG) aims to generate human-like neural questions from an image and potentially other side information (e.g., answer type or the answer itself). Existing works often suffer from the severe one image to many questions mapping problem, which generates uninformative and non-referential questions. Recent work has demonstrated that by leveraging double visual and answer hints, a model can faithfully generate much better quality questions. However, visual hints are not available naturally. Despite they proposed a simple rule-based similarity matching method to obtain candidate visual hints, they could be very noisy practically and thus restrict the quality of generated questions. In this paper, we present a novel learning approach for double-hints based VQG, which can be cast as a weakly supervised learning problem with noises. The key rationale is that the salient visual regions of interest can be viewed as a constraint to improve the generation procedure for producing high-quality questions. As a result, given the predicted salient visual regions of interest, we can focus on estimating the probability of being ground-truth questions, which in turn implicitly measures the quality of predicted visual hints. Experimental results on two benchmark datasets show that our proposed method outperforms the state-of-the-art approaches by a large margin on a variety of metrics, including both automatic machine metrics and human evaluation.

[Pure Exploration in Kernel and Neural Bandits](#)

- Yinglun Zhu, Dongruo Zhou, Ruoxi Jiang, Quanquan Gu, Rebecca Willett, Robert Nowak
- abstract@[open-review](#): We study pure exploration in bandits, where the dimension of the feature representation can be much larger than the number of arms. To overcome the curse of dimensionality, we propose to adaptively embed the feature representation of each arm into a lower-dimensional space and carefully deal with the induced model misspecifications. Our approach is conceptually very different from existing works that can either only handle low-dimensional linear bandits or passively deal with model misspecifications. We showcase the application of our approach to two pure exploration settings that were previously under-studied: (1) the reward function belongs to a possibly infinite-dimensional Reproducing Kernel Hilbert Space, and (2) the reward function is nonlinear and can be approximated by neural networks. Our main results provide sample complexity guarantees that only depend on the effective dimension of the feature spaces in the kernel or neural representations. Extensive experiments conducted on both synthetic and real-world datasets demonstrate the efficacy of our methods.

[Numerical Composition of Differential Privacy](#)

- Sivakanth Gopi, Yin Tat Lee, Lukas Wutschitz
- abstract@[open-review](#): We give a fast algorithm to compose privacy guarantees of differentially private (DP) algorithms to arbitrary accuracy. Our method is based on the notion of privacy loss random variables to quantify the privacy loss of DP algorithms. The running time and memory needed for our algorithm to approximate the privacy curve of a DP algorithm composed with itself $\tilde{O}(k\sqrt{k})$ times is $\tilde{O}(\sqrt{k})$. This improves over the best prior method by Koskela et al. (2020) which requires $\tilde{\Omega}(k^{1.5})$ running time. We demonstrate the utility of our algorithm by accurately computing the privacy loss of DP-SGD algorithm of Abadi et al. (2016) and showing that our algorithm speeds up the privacy computations by a few orders of magnitude compared to prior work, while maintaining similar accuracy.

[Coresets for Classification – Simplified and Strengthened](#)

- Tung Mai, Cameron Musco, Anup Rao
- abstract@[open-review](#): We give relative error coresets for training linear classifiers with a broad class of loss functions, including the logistic loss and hinge loss. Our construction achieves $(1/\epsilon)$ relative error with $\tilde{O}(d \cdot \mu_y(X)^2/\epsilon^2)$ points, where $\mu_y(X)$ is a natural complexity measure of the data matrix $X \in \mathbb{R}^{n \times d}$ and label vector $y \in \{-1, 1\}^n$, introduced by Munteanu et al. 2018. Our result is based on subsampling data points with probabilities proportional to their $\|y\|_1$ Lewis weights. It significantly improves on existing theoretical bounds and performs well in practice, outperforming uniform subsampling along with other importance sampling methods. Our sampling distribution does not depend on the labels, so can be used for active learning. It also does not depend on the specific loss function, so a single coreset can be used in multiple training scenarios.

[Sequential Algorithms for Testing Closeness of Distributions](#)

- Aadil Oufkir, Omar Fawzi, Nicolas Flammarion, Aurélien Garivier
- abstract@[open-review](#): What advantage do sequential procedures provide over batch algorithms for testing properties of unknown distributions? Focusing on the problem of testing whether two distributions D_1 and D_2 on $\{1, \dots, n\}$ are equal or ϵ -far, we give several answers to this question. We show that for a small alphabet size n , there is a sequential algorithm that outperforms any batch algorithm by a factor of at least 4 in terms sample complexity. For a general alphabet size n , we give a sequential algorithm that uses no more samples than its batch counterpart, and possibly fewer if the actual distance between D_1 and D_2 is larger than ϵ . As a corollary, letting ϵ go to 0 , we obtain a sequential algorithm for testing closeness (with no a priori bound on the distance between D_1 and D_2) with a sample complexity $\tilde{O}(\frac{n^{2/3}}{\epsilon^{4/3}} \cdot TV(D_1, D_2)^{4/3})$: this improves over the $\tilde{O}(\frac{n}{\epsilon} \cdot TV(D_1, D_2)^2)$ tester of [Daskalakis and Kawase 2017] and is optimal up to multiplicative constants. We also establish limitations of sequential algorithms for the problem of testing closeness: they can improve the worst case number of samples by at most a constant factor.

[Overlapping Spaces for Compact Graph Representations](#)

- Kirill Shevkunov, Liudmila Prokhorenkova
- abstract@[open-review](#): Various non-trivial spaces are becoming popular for embedding structured data such as graphs, texts, or images. Following spherical and hyperbolic spaces, more general product spaces have been proposed. However, searching for the best configuration of a product space is a resource-intensive procedure, which reduces the practical applicability of the idea. We generalize the concept of product space and introduce an overlapping space that does not have the configuration search problem. The main idea is to allow subsets of coordinates to be shared between spaces of

different types (Euclidean, hyperbolic, spherical). As a result, we often need fewer coordinates to store the objects. Additionally, we propose an optimization algorithm that automatically learns the optimal configuration. Our experiments confirm that overlapping spaces outperform the competitors in graph embedding tasks with different evaluation metrics. We also perform an empirical analysis in a realistic information retrieval setup, where we compare all spaces by incorporating them into DSSM. In this case, the proposed overlapping space consistently achieves nearly optimal results without any configuration tuning. This allows for reducing training time, which can be essential in large-scale applications.

[Hyperparameter Tuning is All You Need for LISTA](#)

- Xiaohan Chen, Jialin Liu, Zhangyang Wang, Wotao Yin
- abstract@[open-review](#): Learned Iterative Shrinkage-Thresholding Algorithm (LISTA) introduces the concept of unrolling an iterative algorithm and training it like a neural network. It has had great success on sparse recovery. In this paper, we show that adding momentum to intermediate variables in the LISTA network achieves a better convergence rate and, in particular, the network with instance-optimal parameters is superlinearly convergent. Moreover, our new theoretical results lead to a practical approach of automatically and adaptively calculating the parameters of a LISTA network layer based on its previous layers. Perhaps most surprisingly, such an adaptive-parameter procedure reduces the training of LISTA to tuning only three hyperparameters from data: a new record set in the context of the recent advances on trimming down LISTA complexity. We call this new ultra-light weight network HyperLISTA. Compared to state-of-the-art LISTA models, HyperLISTA achieves almost the same performance on seen data distributions and performs better when tested on unseen distributions (specifically, those with different sparsity levels and nonzero magnitudes). Code is available: <https://github.com/VITA-Group/HyperLISTA>.

[Foundations of Symbolic Languages for Model Interpretability](#)

- Marcelo Arenas, Daniel Báez, Pablo Barceló, Jorge Pérez, Bernardo Subercaseaux
- abstract@[open-review](#): Several queries and scores have recently been proposed to explain individual predictions over ML models. Examples include queries based on “anchors”, which are parts of an instance that are sufficient to justify its classification, and “feature-perturbation” scores such as SHAP. Given the need for flexible, reliable, and easy-to-apply interpretability methods for ML models, we foresee the need for developing declarative languages to naturally specify different explainability queries. We do this in a principled way by rooting such a language in a logic called FOIL, which allows for expressing many simple but important explainability queries, and might serve as a core for more expressive interpretability languages. We study the computational complexity of FOIL queries over two classes of ML models often deemed to be easily interpretable: decision trees and more general decision diagrams. Since the number of possible inputs for an ML model is exponential in its dimension, tractability of the FOIL evaluation problem is delicate but can be achieved by either restricting the structure of the models, or the fragment of FOIL being evaluated. We also present a prototype implementation of FOIL wrapped in a high-level declarative language and perform experiments showing that such a language can be used in practice.

[Bridging Offline Reinforcement Learning and Imitation Learning: A Tale of Pessimism](#)

- Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, Stuart Russell
- abstract@[open-review](#): Offline (or batch) reinforcement learning (RL) algorithms seek to learn an optimal policy from a fixed dataset without active data collection. Based on the composition of the offline dataset, two main methods are used: imitation learning which is suitable for expert datasets, and vanilla offline RL which often requires uniform coverage datasets. From a practical standpoint, datasets often deviate from these two extremes and the exact data composition is usually unknown. To bridge this gap, we present a new offline RL framework that smoothly interpolates between the two extremes of data composition, hence unifying imitation learning and vanilla offline RL. The new framework is centered around a weak version of the concentrability coefficient that measures the deviation of the behavior policy from the expert policy alone. Under this new framework, we ask: can one develop an algorithm that achieves a minimax optimal rate adaptive to unknown data composition? To address this question, we consider a lower confidence bound (LCB) algorithm developed based on pessimism in the face of uncertainty in offline RL. We study finite-sample properties of LCB as well as information-theoretic limits in multi-armed bandits, contextual bandits, and Markov decision processes (MDPs). Our analysis reveals surprising facts about optimality rates. In particular, in both contextual bandits and RL, LCB achieves a faster rate of $\$1/\sqrt{N}$ for nearly-expert datasets compared to the usual rate of $\$1/\sqrt{N}$ in offline RL, where N is the batch dataset sample size. In contextual bandits with at least two contexts, we prove that LCB is adaptively optimal for the entire data composition range, achieving a smooth transition from imitation learning to offline RL. We further show that LCB is almost adaptively optimal in MDPs.

[Impression learning: Online representation learning with synaptic plasticity](#)

- Colin Bredenberg, Benjamin Lyo, Eero Simoncelli, Cristina Savin
- abstract@[open-review](#): Understanding how the brain constructs statistical models of the sensory world remains a longstanding challenge for computational neuroscience. Here, we derive an unsupervised local synaptic plasticity rule that trains neural circuits to infer latent structure from sensory stimuli via a novel loss function for approximate online Bayesian inference. The learning algorithm is driven by a local error signal computed between two factors that jointly contribute to neural activity: stimulus drive and internal predictions --- the network's 'impression' of the stimulus. Physiologically, we associate these two components with the basal and apical dendrites of pyramidal neurons, respectively. We show that learning can be implemented online, is capable of capturing temporal dependencies in continuous input streams, and generalizes to hierarchical architectures. Furthermore, we demonstrate both analytically and empirically that the algorithm is more data-efficient than a three-factor plasticity alternative, enabling it to learn statistics of high-dimensional, naturalistic inputs. Overall, the model provides a bridge from mechanistic accounts of synaptic plasticity to algorithmic descriptions of unsupervised probabilistic learning and inference.

[How Well do Feature Visualizations Support Causal Understanding of CNN Activations?](#)

- Roland S. Zimmermann, Judy Borowski, Robert Geirhos, Matthias Bethge, Thomas Wallis, Wieland Brendel
- abstract@[open-review](#): A precise understanding of why units in an artificial network respond to certain stimuli would constitute a big step towards explainable artificial intelligence. One widely used approach towards this goal is to visualize unit responses via activation maximization. These feature visualizations are purported to provide humans with precise information about the image features that cause a unit to be activated - an advantage over other alternatives like strongly activating dataset samples. If humans indeed gain causal insight from visualizations, this should enable them to predict the effect of an intervention, such as how occluding a certain patch of the image (say, a dog's head) changes a unit's activation. Here, we test this hypothesis by asking humans to decide which of two square occlusions causes a larger change to a unit's activation. Both a large-scale crowdsourced experiment and measurements with experts show that on average the extremely activating feature visualizations by Olah et al. (2017) indeed help humans on this task ($\$60 \pm 3\%$ accuracy; baseline performance without any visualizations is $\$60 \pm 3\%$). However, they do not provide any substantial advantage over other visualizations (such as e.g. dataset samples), which yield similar performance ($\$66 \pm 3\%$ to $\$67 \pm 3\%$ accuracy). Taken together, we propose an objective psychophysical task to quantify the benefit of unit-level interpretability methods for humans, and find no evidence that a widely-used feature visualization method provides humans with better "causal understanding" of unit activations than simple alternative visualizations.

[Fixes That Fail: Self-Defeating Improvements in Machine-Learning Systems](#)

- Ruihan Wu, Chuan Guo, Awini Hannun, Laurens van der Maaten
- abstract@[open-review](#): Machine-learning systems such as self-driving cars or virtual assistants are composed of a large number of machine-learning models that recognize image content, transcribe speech, analyze natural language, infer preferences, rank options, etc. Models in these systems are often

developed and trained independently, which raises an obvious concern: Can improving a machine-learning model make the overall system worse? We answer this question affirmatively by showing that improving a model can deteriorate the performance of downstream models, even after those downstream models are retrained. Such self-defeating improvements are the result of entanglement between the models in the system. We perform an error decomposition of systems with multiple machine-learning models, which sheds light on the types of errors that can lead to self-defeating improvements. We also present the results of experiments which show that self-defeating improvements emerge in a realistic stereo-based detection system for cars and pedestrians.

[Coarse-to-fine Animal Pose and Shape Estimation](#)

- Chen Li, Gim Hee Lee
- abstract@[open-review](#): Most existing animal pose and shape estimation approaches reconstruct animal meshes with a parametric SMAL model. This is because the low-dimensional pose and shape parameters of the SMAL model makes it easier for deep networks to learn the high-dimensional animal meshes. However, the SMAL model is learned from scans of toy animals with limited pose and shape variations, and thus may not be able to represent highly varying real animals well. This may result in poor fittings of the estimated meshes to the 2D evidences, e.g. 2D keypoints or silhouettes. To mitigate this problem, we propose a coarse-to-fine approach to reconstruct 3D animal mesh from a single image. The coarse estimation stage first estimates the pose, shape and translation parameters of the SMAL model. The estimated meshes are then used as a starting point by a graph convolutional network (GCN) to predict a per-vertex deformation in the refinement stage. This combination of SMAL-based and vertex-based representations benefits from both parametric and non-parametric representations. We design our mesh refinement GCN (MRGCN) as an encoder-decoder structure with hierarchical feature representations to overcome the limited receptive field of traditional GCNs. Moreover, we observe that the global image feature used by existing animal mesh reconstruction works is unable to capture detailed shape information for mesh refinement. We thus introduce a local feature extractor to retrieve a vertex-level feature and use it together with the global feature as the input of the MRGCN. We test our approach on the StanfordExtra dataset and achieve state-of-the-art results. Furthermore, we test the generalization capacity of our approach on the Animal Pose and BADJA datasets. Our code is available at the project website.

[Meta-Learning Sparse Implicit Neural Representations](#)

- Jaeho Lee, Jihoon Tack, Namhoon Lee, Jinwoo Shin
- abstract@[open-review](#): Implicit neural representations are a promising new avenue of representing general signals by learning a continuous function that, parameterized as a neural network, maps the domain of a signal to its codomain; the mapping from spatial coordinates of an image to its pixel values, for example. Being capable of conveying fine details in a high dimensional signal, unboundedly of its domain, implicit neural representations ensure many advantages over conventional discrete representations. However, the current approach is difficult to scale for a large number of signals or a data set, since learning a neural representation---which is parameter heavy by itself---for each signal individually requires a lot of memory and computations. To address this issue, we propose to leverage a meta-learning approach in combination with network compression under a sparsity constraint, such that it renders a well-initialized sparse parameterization that evolves quickly to represent a set of unseen signals in the subsequent training. We empirically demonstrate that meta-learned sparse neural representations achieve a much smaller loss than dense meta-learned models with the same number of parameters, when trained to fit each signal using the same number of optimization steps.

[Rethinking Space-Time Networks with Improved Memory Coverage for Efficient Video Object Segmentation](#)

- Ho Kei Cheng, Yu-Wing Tai, Chi-Keung Tang
- abstract@[open-review](#): This paper presents a simple yet effective approach to modeling space-time correspondences in the context of video object segmentation. Unlike most existing approaches, we establish correspondences directly between frames without re-encoding the mask features for every object, leading to a highly efficient and robust framework. With the correspondences, every node in the current query frame is inferred by aggregating features from the past in an associative fashion. We cast the aggregation process as a voting problem and find that the existing inner-product affinity leads to poor use of memory with a small (fixed) subset of memory nodes dominating the votes, regardless of the query. In light of this phenomenon, we propose using the negative squared Euclidean distance instead to compute the affinities. We validated that every memory node now has a chance to contribute, and experimentally showed that such diversified voting is beneficial to both memory efficiency and inference accuracy. The synergy of correspondence networks and diversified voting works exceedingly well, achieves new state-of-the-art results on both DAVIS and YouTubeVOS datasets while running significantly faster at 20+ FPS for multiple objects without bells and whistles.

[Sparse Spiking Gradient Descent](#)

- Nicolas Perez-Nieves, Dan Goodman
- abstract@[open-review](#): There is an increasing interest in emulating Spiking Neural Networks (SNNs) on neuromorphic computing devices due to their low energy consumption. Recent advances have allowed training SNNs to a point where they start to compete with traditional Artificial Neural Networks (ANNs) in terms of accuracy, while at the same time being energy efficient when run on neuromorphic hardware. However, the process of training SNNs is still based on dense tensor operations originally developed for ANNs which do not leverage the spatiotemporally sparse nature of SNNs. We present here the first sparse SNN backpropagation algorithm which achieves the same or better accuracy as current state of the art methods while being significantly faster and more memory efficient. We show the effectiveness of our method on real datasets of varying complexity (Fashion-MNIST, Neuromorphic-MNIST and Spiking Heidelberg Digits) achieving a speedup in the backward pass of up to \$150\times\$, and \$85\%\$ more memory efficient, without losing accuracy.

[Rethinking Calibration of Deep Neural Networks: Do Not Be Afraid of Overconfidence](#)

- Deng-Bao Wang, Lei Feng, Min-Ling Zhang
- abstract@[open-review](#): Capturing accurate uncertainty quantification of the prediction from deep neural networks is important in many real-world decision-making applications. A reliable predictor is expected to be accurate when it is confident about its predictions and indicate high uncertainty when it is likely to be inaccurate. However, modern neural networks have been found to be poorly calibrated, primarily in the direction of overconfidence. In recent years, there is a surge of research on model calibration by leveraging implicit or explicit regularization techniques during training, which obtain well calibration by avoiding overconfident outputs. In our study, we empirically found that despite the predictions obtained from these regularized models are better calibrated, they suffer from not being as calibratable, namely, it is harder to further calibrate their predictions with post-hoc calibration methods like temperature scaling and histogram binning. We conduct a series of empirical studies showing that overconfidence may not hurt final calibration performance if post-hoc calibration is allowed, rather, the penalty of confident outputs will compress the room of potential improvements in post-hoc calibration phase. Our experimental findings point out a new direction to improve calibration of DNNs by considering main training and post-hoc calibration as a unified framework.

[Towards Efficient and Effective Adversarial Training](#)

- Gaurang Sriramanan, Sravanti Addepalli, Arya Baburaj, Venkatesh Babu R
- abstract@[open-review](#): The vulnerability of Deep Neural Networks to adversarial attacks has spurred immense interest towards improving their robustness. However, present state-of-the-art adversarial defenses involve the use of 10-step adversaries during training, which renders them computationally infeasible for application to large-scale datasets. While the recent single-step defenses show promising direction, their robustness is not

on par with multi-step training methods. In this work, we bridge this performance gap by introducing a novel Nuclear-Norm regularizer on network predictions to enforce function smoothing in the vicinity of data samples. While prior works consider each data sample independently, the proposed regularizer uses the joint statistics of adversarial samples across a training minibatch to enhance optimization during both attack generation and training, obtaining state-of-the-art results amongst efficient defenses. We achieve further gains by incorporating exponential averaging of network weights over training iterations. We finally introduce a Hybrid training approach that combines the effectiveness of a two-step variant of the proposed defense with the efficiency of a single-step defense. We demonstrate superior results when compared to multi-step defenses such as TRADES and PGD-AT as well, at a significantly lower computational cost.

[Intriguing Properties of Contrastive Losses](#)

- Ting Chen, Calvin Luo, Lala Li
- abstract@[open-review](#): We study three intriguing properties of contrastive learning. First, we generalize the standard contrastive loss to a broader family of losses, and we find that various instantiations of the generalized loss perform similarly under the presence of a multi-layer non-linear projection head. Second, we study if instance-based contrastive learning (with a global image representation) can learn well on images with multiple objects present. We find that meaningful hierarchical local features can be learned despite the fact that these objectives operate on global instance-level features. Finally, we study the phenomenon of feature suppression among competing features shared across augmented views, such as "color distribution" vs "object class". We construct datasets with explicit and controllable competing features, and show that, for contrastive learning, a few bits of easy-to-learn shared features can suppress, and even fully prevent, the learning of other sets of competing features. In scenarios where there are multiple objects in an image, the dominant object would suppress the learning of smaller objects. Existing contrastive learning methods critically rely on data augmentation to favor certain sets of features over others, and could suffer from learning saturation for scenarios where existing augmentations cannot fully address the feature suppression. This poses open challenges to existing contrastive learning techniques.

[Detecting Moments and Highlights in Videos via Natural Language Queries](#)

- Jie Lei, Tamara L Berg, Mohit Bansal
- abstract@[open-review](#): Detecting customized moments and highlights from videos given natural language (NL) user queries is an important but under-studied topic. One of the challenges in pursuing this direction is the lack of annotated data. To address this issue, we present the Query-based Video Highlights (QVHighlights) dataset. It consists of over 10,000 YouTube videos, covering a wide range of topics, from everyday activities and travel in lifestyle vlog videos to social and political activities in news videos. Each video in the dataset is annotated with: (1) a human-written free-form NL query, (2) relevant moments in the video w.r.t. the query, and (3) five-point scale saliency scores for all query-relevant clips. This comprehensive annotation enables us to develop and evaluate systems that detect relevant moments as well as salient highlights for diverse, flexible user queries. We also present a strong baseline for this task, Moment-DETR, a transformer encoder-decoder model that views moment retrieval as a direct set prediction problem, taking extracted video and query representations as inputs and predicting moment coordinates and saliency scores end-to-end. While our model does not utilize any human prior, we show that it performs competitively when compared to well-engineered architectures. With weakly supervised pretraining using ASR captions, Moment-DETR substantially outperforms previous methods. Lastly, we present several ablations and visualizations of Moment-DETR. Data and code is publicly available at https://github.com/jayleicn/moment_detr.

[Stochastic optimization under time drift: iterate averaging, step-decay schedules, and high probability guarantees](#)

- Joshua Cutler, Dmitriy Drusvyatskiy, Zaid Harchaoui
- abstract@[open-review](#): We consider the problem of minimizing a convex function that is evolving in time according to unknown and possibly stochastic dynamics. Such problems abound in the machine learning and signal processing literature, under the names of concept drift and stochastic tracking. We provide novel non-asymptotic convergence guarantees for stochastic algorithms with iterate averaging, focusing on bounds valid both in expectation and with high probability. Notably, we show that the tracking efficiency of the proximal stochastic gradient method depends only logarithmically on the initialization quality when equipped with a step-decay schedule.

[Learning Stable Deep Dynamics Models for Partially Observed or Delayed Dynamical Systems](#)

- Andreas Schlaginhaufen, Philippe Wenk, Andreas Krause, Florian Dorfler
- abstract@[open-review](#): Learning how complex dynamical systems evolve over time is a key challenge in system identification. For safety critical systems, it is often crucial that the learned model is guaranteed to converge to some equilibrium point. To this end, neural ODEs regularized with neural Lyapunov functions are a promising approach when states are fully observed. For practical applications however, {\em partial observations} are the norm. As we will demonstrate, initialization of unobserved augmented states can become a key problem for neural ODEs. To alleviate this issue, we propose to augment the system's state with its history. Inspired by state augmentation in discrete-time systems, we thus obtain {\em neural delay differential equations}. Based on classical time delay stability analysis, we then show how to ensure stability of the learned models, and theoretically analyze our approach. Our experiments demonstrate its applicability to stable system identification of partially observed systems and learning a stabilizing feedback policy in delayed feedback control.

[An Uncertainty Principle is a Price of Privacy-Preserving Microdata](#)

- John Abowd, Robert Ashmead, Ryan Cumings-Menon, Simson Garfinkel, Daniel Kifer, Philip Leclerc, William Sexton, Ashley Simpson, Christine Task, Pavel Zhuravlev
- abstract@[open-review](#): Privacy-protected microdata are often the desired output of a differentially private algorithm since microdata is familiar and convenient for downstream users. However, there is a statistical price for this kind of convenience. We show that an uncertainty principle governs the trade-off between accuracy for a population of interest (sum query'') vs. accuracy for its component sub-populations (point queries"). Compared to differentially private query answering systems that are not required to produce microdata, accuracy can degrade by a logarithmic factor. For example, in the case of pure differential privacy, without the microdata requirement, one can provide noisy answers to the sum query and all point queries while guaranteeing that each answer has squared error $\$O(1/\epsilon^2)$. With the microdata requirement, one must choose between allowing an additional $\$log^2(d)$ factor (d is the number of point queries) for some point queries or allowing an extra $\$O(d^2)$ factor for the sum query. We present lower bounds for pure, approximate, and concentrated differential privacy. We propose mitigation strategies and create a collection of benchmark datasets that can be used for public study of this problem.

[Fairness in Ranking under Uncertainty](#)

- Ashudeep Singh, David Kempe, Thorsten Joachims
- abstract@[open-review](#): Fairness has emerged as an important consideration in algorithmic decision making. Unfairness occurs when an agent with higher merit obtains a worse outcome than an agent with lower merit. Our central point is that a primary cause of unfairness is uncertainty. A principal or algorithm making decisions never has access to the agents' true merit, and instead uses proxy features that only imperfectly predict merit (e.g., GPA, star ratings, recommendation letters). None of these ever fully capture an agent's merit; yet existing approaches have mostly been defining fairness notions directly based on observed features and outcomes. Our primary point is that it is more principled to acknowledge and model the uncertainty explicitly. The role of observed features is to give rise to a posterior distribution of the agents' merits. We use this viewpoint to define a notion of approximate fairness in ranking. We call an algorithm $\$phi$-fair$ (for $\$phi \in [0,1]$) if it has the following property for all agents x and all k : if agent x is among the top

\$k\$ agents with respect to merit with probability at least \$\\rho\$ (according to the posterior merit distribution), then the algorithm places the agent among the top \$k\$ agents in its ranking with probability at least \$\\phi \\rho\$. We show how to compute rankings that optimally trade off approximate fairness against utility to the principal. In addition to the theoretical characterization, we present an empirical analysis of the potential impact of the approach in simulation studies. For real-world validation, we applied the approach in the context of a paper recommendation system that we built and fielded at the KDD 2020 conference.

[Generalized Proximal Policy Optimization with Sample Reuse](#)

- James Queeney, Yannis Paschalidis, Christos G Cassandras
- abstract@[open-review](#): In real-world decision making tasks, it is critical for data-driven reinforcement learning methods to be both stable and sample efficient. On-policy methods typically generate reliable policy improvement throughout training, while off-policy methods make more efficient use of data through sample reuse. In this work, we combine the theoretically supported stability benefits of on-policy algorithms with the sample efficiency of off-policy algorithms. We develop policy improvement guarantees that are suitable for the off-policy setting, and connect these bounds to the clipping mechanism used in Proximal Policy Optimization. This motivates an off-policy version of the popular algorithm that we call Generalized Proximal Policy Optimization with Sample Reuse. We demonstrate both theoretically and empirically that our algorithm delivers improved performance by effectively balancing the competing goals of stability and sample efficiency.

[Mosaicking to Distill: Knowledge Distillation from Out-of-Domain Data](#)

- Gongfan Fang, Yifan Bao, Jie Song, Xinchao Wang, Donglin Xie, Chengchao Shen, Mingli Song
- abstract@[open-review](#): Knowledge distillation~(KD) aims to craft a compact student model that imitates the behavior of a pre-trained teacher in a target domain. Prior KD approaches, despite their gratifying results, have largely relied on the premise that \emph{in-domain} data is available to carry out the knowledge transfer. Such an assumption, unfortunately, in many cases violates the practical setting, since the original training data or even the data domain is often unreachable due to privacy or copyright reasons. In this paper, we attempt to tackle an ambitious task, termed as \emph{out-of-domain} knowledge distillation~(OOD-KD), which allows us to conduct KD using only OOD data that can be readily obtained at a very low cost. Admittedly, OOD-KD is by nature a highly challenging task due to the agnostic domain gap. To this end, we introduce a handy yet surprisingly efficacious approach, dubbed as \textit{MosaicKD}. The key insight behind MosaicKD lies in that, samples from various domains share common local patterns, even though their global semantic may vary significantly; these shared local patterns, in turn, can be re-assembled analogous to mosaic tiling, to approximate the in-domain data and to further alleviating the domain discrepancy. In MosaicKD, this is achieved through a four-player min-max game, in which a generator, a discriminator, a student network, are collectively trained in an adversarial manner, partially under the guidance of a pre-trained teacher. We validate MosaicKD over \{classification and semantic segmentation tasks\} across various benchmarks, and demonstrate that it yields results much superior to the state-of-the-art counterparts on OOD data. Our code is available at \url{https://github.com/zju-vipa/MosaicKD}.

[Batch Active Learning at Scale](#)

- Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Afshin Rostamizadeh, Sanjiv Kumar
- abstract@[open-review](#): The ability to train complex and highly effective models often requires an abundance of training data, which can easily become a bottleneck in cost, time, and computational resources. Batch active learning, which adaptively issues batched queries to a labeling oracle, is a common approach for addressing this problem. The practical benefits of batch sampling come with the downside of less adaptivity and the risk of sampling redundant examples within a batch -- a risk that grows with the batch size. In this work, we analyze an efficient active learning algorithm, which focuses on the large batch setting. In particular, we show that our sampling method, which combines notions of uncertainty and diversity, easily scales to batch sizes (100K-1M) several orders of magnitude larger than used in previous studies and provides significant improvements in model training efficiency compared to recent baselines. Finally, we provide an initial theoretical analysis, proving label complexity guarantees for a related sampling method, which we show is approximately equivalent to our sampling method in specific settings.

[Joint Semantic Mining for Weakly Supervised RGB-D Salient Object Detection](#)

- Jingjing Li, Wei Ji, Qi Bi, Cheng Yan, Miao Zhang, Yongri Piao, Huchuan Lu, Li cheng
- abstract@[open-review](#): Training saliency detection models with weak supervisions, e.g., image-level tags or captions, is appealing as it removes the costly demand of per-pixel annotations. Despite the rapid progress of RGB-D saliency detection in fully-supervised setting, it however remains an unexplored territory when only weak supervision signals are available. This paper is set to tackle the problem of weakly-supervised RGB-D salient object detection. The key insight in this effort is the idea of maintaining per-pixel pseudo-labels with iterative refinements by reconciling the multimodal input signals in our joint semantic mining (JSM). Considering the large variations in the raw depth map and the lack of explicit pixel-level supervisions, we propose spatial semantic modeling (SSM) to capture saliency-specific depth cues from the raw depth and produce depth-refined pseudo-labels. Moreover, tags and captions are incorporated via a fill-in-the-blank training in our textual semantic modeling (TSM) to estimate the confidences of competing pseudo-labels. At test time, our model involves only a light-weight sub-network of the training pipeline, i.e., it requires only an RGB image as input, thus allowing efficient inference. Extensive evaluations demonstrate the effectiveness of our approach under the weakly-supervised setting. Importantly, our method could also be adapted to work in both fully-supervised and unsupervised paradigms. In each of these scenarios, superior performance has been attained by our approach with comparing to the state-of-the-art dedicated methods. As a by-product, a CapS dataset is constructed by augmenting existing benchmark training set with additional image tags and captions.

[Not All Images are Worth 16x16 Words: Dynamic Transformers for Efficient Image Recognition](#)

- Yulin Wang, Rui Huang, Shiji Song, Zeyi Huang, Gao Huang
- abstract@[open-review](#): Vision Transformers (ViT) have achieved remarkable success in large-scale image recognition. They split every 2D image into a fixed number of patches, each of which is treated as a token. Generally, representing an image with more tokens would lead to higher prediction accuracy, while it also results in drastically increased computational cost. To achieve a decent trade-off between accuracy and speed, the number of tokens is empirically set to 16x16 or 14x14. In this paper, we argue that every image has its own characteristics, and ideally the token number should be conditioned on each individual input. In fact, we have observed that there exist a considerable number of “easy” images which can be accurately predicted with a mere number of 4x4 tokens, while only a small fraction of “hard” ones need a finer representation. Inspired by this phenomenon, we propose a Dynamic Transformer to automatically configure a proper number of tokens for each input image. This is achieved by cascading multiple Transformers with increasing numbers of tokens, which are sequentially activated in an adaptive fashion at test time, i.e., the inference is terminated once a sufficiently confident prediction is produced. We further design efficient feature reuse and relationship reuse mechanisms across different components of the Dynamic Transformer to reduce redundant computations. Extensive empirical results on ImageNet, CIFAR-10, and CIFAR-100 demonstrate that our method significantly outperforms the competitive baselines in terms of both theoretical computational efficiency and practical inference speed. Code and pre-trained models (based on PyTorch and MindSpore) are available at <https://github.com/blackfeather-wang/Dynamic-Vision-Transformer> and <https://github.com/blackfeather-wang/Dynamic-Vision-Transformer-MindSpore>.

[Contrastive Learning for Neural Topic Model](#)

- Thong Nguyen, Anh Tuan Luu

- abstract@[open-review](#): Recent empirical studies show that adversarial topic models (ATM) can successfully capture semantic patterns of the document by differentiating a document with another dissimilar sample. However, utilizing that discriminative-generative architecture has two important drawbacks: (1) the architecture does not relate similar documents, which has the same document-word distribution of salient words; (2) it restricts the ability to integrate external information, such as sentiments of the document, which has been shown to benefit the training of neural topic model. To address those issues, we revisit the adversarial topic architecture in the view point of mathematical analysis, propose a novel approach to re-formulate discriminative goal as an optimization problem, and design a novel sampling method which facilitates the integration of external variables. The reformulation encourages the model to incorporate the relations among similar samples and enforces the constraint on the similarity among dissimilar ones; while the sampling method, which is based on the internal input and reconstructed output, helps inform the model of salient words contributing to the main topic. Experimental results show that our framework outperforms other state-of-the-art neural topic models in three common benchmark datasets that belong to various domains, vocabulary sizes, and document lengths in terms of topic coherence.

[Learning in two-player zero-sum partially observable Markov games with perfect recall](#)

- Tadashi Kozuno, Pierre Ménard, Remi Munos, Michal Valko
- abstract@[open-review](#): We study the problem of learning a Nash equilibrium (NE) in an extensive game with imperfect information (EGII) through self-play. Precisely, we focus on two-player, zero-sum, episodic, tabular EGII under the \textit{perfect-recall} assumption where the only feedback is realizations of the game (bandit feedback). In particular the \textit{dynamics of the EGII is not known}---we can only access it by sampling or interacting with a game simulator. For this learning setting, we provide the Implicit Exploration Online Mirror Descent (IXOMD) algorithm. It is a model-free algorithm with a high-probability bound on convergence rate to the NE of order $\$1/\sqrt{T}$ where $\sim T$ is the number of played games. Moreover IXOMD is computationally efficient as it needs to perform the updates only along the sampled trajectory.

[A Geometric Structure of Acceleration and Its Role in Making Gradients Small Fast](#)

- Jongmin Lee, Chanwoo Park, Ernest Ryu
- abstract@[open-review](#): Since Nesterov's seminal 1983 work, many accelerated first-order optimization methods have been proposed, but their analyses lacks a common unifying structure. In this work, we identify a geometric structure satisfied by a wide range of first-order accelerated methods. Using this geometric insight, we present several novel generalizations of accelerated methods. Most interesting among them is a method that reduces the squared gradient norm with $\mathcal{O}(1/K^4)$ rate in the prox-grad setup, faster than the $\mathcal{O}(1/K^3)$ rates of Nesterov's FGM or Kim and Fessler's FPGM-m.

[ATISS: Autoregressive Transformers for Indoor Scene Synthesis](#)

- Despoina Paschalidou, Amlan Kar, Maria Shugrina, Karsten Kreis, Andreas Geiger, Sanja Fidler
- abstract@[open-review](#): The ability to synthesize realistic and diverse indoor furniture layouts automatically or based on partial input, unlocks many applications, from better interactive 3D tools to data synthesis for training and simulation. In this paper, we present ATISS, a novel autoregressive transformer architecture for creating diverse and plausible synthetic indoor environments, given only the room type and its floor plan. In contrast to prior work, which poses scene synthesis as sequence generation, our model generates rooms as unordered sets of objects. We argue that this formulation is more natural, as it makes ATISS generally useful beyond fully automatic room layout synthesis. For example, the same trained model can be used in interactive applications for general scene completion, partial room re-arrangement with any objects specified by the user, as well as object suggestions for any partial room. To enable this, our model leverages the permutation equivariance of the transformer when conditioning on the partial scene, and is trained to be permutation-invariant across object orderings. Our model is trained end-to-end as an autoregressive generative model using only labeled 3D bounding boxes as supervision. Evaluations on four room types in the 3D-FRONT dataset demonstrate that our model consistently generates plausible room layouts that are more realistic than existing methods. In addition, it has fewer parameters, is simpler to implement and train and runs up to 8 times faster than existing methods.

[Generalized Depthwise-Separable Convolutions for Adversarially Robust and Efficient Neural Networks](#)

- Hassan Dbouk, Naresh Shanbhag
- abstract@[open-review](#): Despite their tremendous successes, convolutional neural networks (CNNs) incur high computational/storage costs and are vulnerable to adversarial perturbations. Recent works on robust model compression address these challenges by combining model compression techniques with adversarial training. But these methods are unable to improve throughput (frames-per-second) on real-life hardware while simultaneously preserving robustness to adversarial perturbations. To overcome this problem, we propose the method of Generalized Depthwise-Separable (GDWS) convolution - an efficient, universal, post-training approximation of a standard 2D convolution. GDWS dramatically improves the throughput of a standard pre-trained network on real-life hardware while preserving its robustness. Lastly, GDWS is scalable to large problem sizes since it operates on pre-trained models and doesn't require any additional training. We establish the optimality of GDWS as a 2D convolution approximator and present exact algorithms for constructing optimal GDWS convolutions under complexity and error constraints. We demonstrate the effectiveness of GDWS via extensive experiments on CIFAR-10, SVHN, and ImageNet datasets. Our code can be found at <https://github.com/hsndbk4/GDWS>.

[A Provably Efficient Model-Free Posterior Sampling Method for Episodic Reinforcement Learning](#)

- Christoph Dann, Mehryar Mohri, Tong Zhang, Julian Zimmert
- abstract@[open-review](#): Thompson Sampling is one of the most effective methods for contextual bandits and has been generalized to posterior sampling for certain MDP settings. However, existing posterior sampling methods for reinforcement learning are limited by being model-based or lack worst-case theoretical guarantees beyond linear MDPs. This paper proposes a new model-free formulation of posterior sampling that applies to more general episodic reinforcement learning problems with theoretical guarantees. We introduce novel proof techniques to show that under suitable conditions, the worst-case regret of our posterior sampling method matches the best known results of optimization based methods. In the linear MDP setting with dimension, the regret of our algorithm scales linearly with the dimension as compared to a quadratic dependence of the existing posterior sampling-based exploration algorithms.

[Fast Federated Learning in the Presence of Arbitrary Device Unavailability](#)

- Xinran Gu, Kaixuan Huang, Jingzhao Zhang, Longbo Huang
- abstract@[open-review](#): Federated learning (FL) coordinates with numerous heterogeneous devices to collaboratively train a shared model while preserving user privacy. Despite its multiple advantages, FL faces new challenges. One challenge arises when devices drop out of the training process. In this case, the convergence of popular FL algorithms such as FedAvg is severely influenced by the straggling devices. To tackle this challenge, we study federated learning algorithms in the presence of arbitrary device unavailability and propose an algorithm named Memory-augmented Impatient Federated Averaging (MIFA). Our algorithm efficiently avoids excessive latency induced by inactive devices, and corrects the gradient bias using the memorized latest updates from them. We prove that MIFA achieves minimax optimal convergence rates on non-i.i.d. data for both strongly convex and non-convex smooth functions. We also provide an explicit characterization of the improvement over baseline algorithms through a case study, and validate the results by numerical experiments on real-world datasets.

[On The Structure of Parametric Tournaments with Application to Ranking from Pairwise Comparisons](#)

- Vishnu Veerathu, Arun Rajkumar
- abstract@[open-review](#): We consider the classical problem of finding the minimum feedback arc set on tournaments (MFAST). The problem is NP-hard in general and we study it for important classes of tournaments that arise naturally in the problem of learning to rank from pairwise comparisons. Specifically, we consider tournaments classes that arise out of parametric preference matrices that can lead to cyclic preference relations. We investigate their structural properties via forbidden sub tournament configurations. Towards this, we introduce \emph{Tournament Dimension} - a combinatorial parameter that characterizes the size of a forbidden configuration for rank \$r\$ tournament classes i.e., classes that arise out pairwise preference matrices which lead to rank \$r\$ skew-symmetric matrices under a suitable link function. Our main result is a polynomial-time algorithm - \texttt{Rank2Rank} - that solves the MFAST problem for the rank \$2\$ tournament class. This is achieved via a geometric characterization that relies on our explicit construction of a forbidden configuration for this class. Building on our understanding of the rank-\$2\$ tournament class, we propose a very general and flexible parametric pairwise preference model called the local-global model which subsumes the popular Bradley-Terry-Luce/Thurstone classes to capture locally cyclic as well as globally acyclic preference relations. We develop a polynomial-time algorithm - \texttt{BlockRank2Rank} - to solve the MFAST problem on the associated Block-Rank \$2\$ tournament class. As an application, we study the problem of learning to rank from pairwise comparisons under the proposed local-global preference model. Exploiting our structural characterization, we propose \texttt{PairwiseBlockRank} - a pairwise ranking algorithm for this class. We show sample complexity bounds of \texttt{PairwiseBlockRank} to learn a good ranking under the proposed model. Finally, we conduct experiments on synthetic and real-world datasets to show the efficacy of the proposed algorithm.

[SegFormer: Simple and Efficient Design for Semantic Segmentation with Transformers](#)

- Enze Xie, Wenhui Wang, Zhiding Yu, Anima Anandkumar, Jose M. Alvarez, Ping Luo
- abstract@[open-review](#): We present SegFormer, a simple, efficient yet powerful semantic segmentation framework which unifies Transformers with lightweight multilayer perceptron (MLP) decoders. SegFormer has two appealing features: 1) SegFormer comprises a novel hierarchically structured Transformer encoder which outputs multiscale features. It does not need positional encoding, thereby avoiding the interpolation of positional codes which leads to decreased performance when the testing resolution differs from training. 2) SegFormer avoids complex decoders. The proposed MLP decoder aggregates information from different layers, and thus combining both local attention and global attention to render powerful representations. We show that this simple and lightweight design is the key to efficient segmentation on Transformers. We scale our approach up to obtain a series of models from SegFormer-B0 to Segformer-B5, which reaches much better performance and efficiency than previous counterparts. For example, SegFormer-B4 achieves 50.3% mIoU on ADE20K with 64M parameters, being 5x smaller and 2.2% better than the previous best method. Our best model, SegFormer-B5, achieves 84.0% mIoU on Cityscapes validation set and shows excellent zero-shot robustness on Cityscapes-C.

[Fairness via Representation Neutralization](#)

- Mengnan Du, Subhabrata Mukherjee, Guanchu Wang, Ruixiang Tang, Ahmed Awadallah, Xia Hu
- abstract@[open-review](#): Existing bias mitigation methods for DNN models primarily work on learning debiased encoders. This process not only requires a lot of instance-level annotations for sensitive attributes, it also does not guarantee that all fairness sensitive information has been removed from the encoder. To address these limitations, we explore the following research question: Can we reduce the discrimination of DNN models by only debiasing the classification head, even with biased representations as inputs? To this end, we propose a new mitigation technique, namely, Representation Neutralization for Fairness (RNF) that achieves fairness by debiasing only the task-specific classification head of DNN models. To this end, we leverage samples with the same ground-truth label but different sensitive attributes, and use their neutralized representations to train the classification head of the DNN model. The key idea of RNF is to discourage the classification head from capturing spurious correlation between fairness sensitive information in encoder representations with specific class labels. To address low-resource settings with no access to sensitive attribute annotations, we leverage a bias-amplified model to generate proxy annotations for sensitive attributes. Experimental results over several benchmark datasets demonstrate our RNF framework to effectively reduce discrimination of DNN models with minimal degradation in task-specific performance.

[Residual Relaxation for Multi-view Representation Learning](#)

- Yifei Wang, Zhengyang Geng, Feng Jiang, Chuming Li, Yisen Wang, Jiansheng Yang, Zhouchen Lin
- abstract@[open-review](#): Multi-view methods learn representations by aligning multiple views of the same image and their performance largely depends on the choice of data augmentation. In this paper, we notice that some other useful augmentations, such as image rotation, are harmful for multi-view methods because they cause a semantic shift that is too large to be aligned well. This observation motivates us to relax the exact alignment objective to better cultivate stronger augmentations. Taking image rotation as a case study, we develop a generic approach, Pretext-aware Residual Relaxation (Prelax), that relaxes the exact alignment by allowing an adaptive residual vector between different views and encoding the semantic shift through pretext-aware learning. Extensive experiments on different backbones show that our method can not only improve multi-view methods with existing augmentations, but also benefit from stronger image augmentations like rotation.

[Do Vision Transformers See Like Convolutional Neural Networks?](#)

- Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, Alexey Dosovitskiy
- abstract@[open-review](#): Convolutional neural networks (CNNs) have so far been the de-facto model for visual data. Recent work has shown that (Vision) Transformer models (ViT) can achieve comparable or even superior performance on image classification tasks. This raises a central question: how are Vision Transformers solving these tasks? Are they acting like convolutional networks, or learning entirely different visual representations? Analyzing the internal representation structure of ViTs and CNNs on image classification benchmarks, we find striking differences between the two architectures, such as ViT having more uniform representations across all layers. We explore how these differences arise, finding crucial roles played by self-attention, which enables early aggregation of global information, and ViT residual connections, which strongly propagate features from lower to higher layers. We study the ramifications for spatial localization, demonstrating ViTs successfully preserve input spatial information, with noticeable effects from different classification methods. Finally, we study the effect of (pretraining) dataset scale on intermediate features and transfer learning, and conclude with a discussion on connections to new architectures such as the MLP-Mixer.

[Optimization-Based Algebraic Multigrid Coarsening Using Reinforcement Learning](#)

- Ali Taghibakhshi, Scott MacLachlan, Luke Olson, Matthew West
- abstract@[open-review](#): Large sparse linear systems of equations are ubiquitous in science and engineering, such as those arising from discretizations of partial differential equations. Algebraic multigrid (AMG) methods are one of the most common methods of solving such linear systems, with an extensive body of underlying mathematical theory. A system of linear equations defines a graph on the set of unknowns and each level of a multigrid solver requires the selection of an appropriate coarse graph along with restriction and interpolation operators that map to and from the coarse representation. The efficiency of the multigrid solver depends critically on this selection and many selection methods have been developed over the years. Recently, it has been demonstrated that it is possible to directly learn the AMG interpolation and restriction operators, given a coarse graph selection. In this paper, we consider the complementary problem of learning to coarsen graphs for a multigrid solver, a necessary step in developing fully learnable AMG methods. We propose a method using a reinforcement learning (RL) agent based on graph neural networks (GNNs), which can learn to perform graph coarsening on small planar training graphs and then be applied to unstructured large planar graphs, assuming bounded node degree. We demonstrate that this method can

produce better coarse graphs than existing algorithms, even as the graph size increases and other properties of the graph are varied. We also propose an efficient inference procedure for performing graph coarsening that results in linear time complexity in graph size.

[Delayed Propagation Transformer: A Universal Computation Engine towards Practical Control in Cyber-Physical Systems](#)

- Wenqing Zheng, Qiangqiang Guo, Hao Yang, Peihao Wang, Zhangyang Wang
- abstract@[open-review](#): Multi-agent control is a central theme in the Cyber-Physical Systems (CPS). However, current control methods either receive non-Markovian states due to insufficient sensing and decentralized design, or suffer from poor convergence. This paper presents the Delayed Propagation Transformer (DePT), a new transformer-based model that specializes in the global modeling of CPS while taking into account the immutable constraints from the physical world. DePT induces a cone-shaped spatial-temporal attention prior, which injects the information propagation and aggregation principles and enables a global view. With physical constraint inductive bias baked into its design, our DePT is ready to plug and play for a broad class of multi-agent systems. The experimental results on one of the most challenging CPS -- network-scale traffic signal control system in the open world -- show that our model outperformed the state-of-the-art expert methods on synthetic and real-world datasets. Our codes are released at: <https://github.com/VITA-Group/DePT>.

[Explaining Latent Representations with a Corpus of Examples](#)

- Jonathan Crabbe, Zhaozhi Qian, Fergus Imrie, Mihaela van der Schaar
- abstract@[open-review](#): Modern machine learning models are complicated. Most of them rely on convoluted latent representations of their input to issue a prediction. To achieve greater transparency than a black-box that connects inputs to predictions, it is necessary to gain a deeper understanding of these latent representations. To that aim, we propose SimplEx: a user-centred method that provides example-based explanations with reference to a freely selected set of examples, called the corpus. SimplEx uses the corpus to improve the user's understanding of the latent space with post-hoc explanations answering two questions: (1) Which corpus examples explain the prediction issued for a given test example? (2) What features of these corpus examples are relevant for the model to relate them to the test example? SimplEx provides an answer by reconstructing the test latent representation as a mixture of corpus latent representations. Further, we propose a novel approach, the integrated Jacobian, that allows SimplEx to make explicit the contribution of each corpus feature in the mixture. Through experiments on tasks ranging from mortality prediction to image classification, we demonstrate that these decompositions are robust and accurate. With illustrative use cases in medicine, we show that SimplEx empowers the user by highlighting relevant patterns in the corpus that explain model representations. Moreover, we demonstrate how the freedom in choosing the corpus allows the user to have personalized explanations in terms of examples that are meaningful for them.

[Explaining heterogeneity in medial entorhinal cortex with task-driven neural networks](#)

- Aran Nayebi, Alexander Attinger, Malcolm Campbell, Kiah Hardcastle, Isabel Low, Caitlin S Mallory, Gabriel Mel, Ben Sorscher, Alex H Williams, Surya Ganguli, Lisa Giocomo, Dan Yamins
- abstract@[open-review](#): Medial entorhinal cortex (MEC) supports a wide range of navigational and memory related behaviors. Well-known experimental results have revealed specialized cell types in MEC --- e.g. grid, border, and head-direction cells --- whose highly stereotypical response profiles are suggestive of the role they might play in supporting MEC functionality. However, the majority of MEC neurons do not exhibit stereotypical firing patterns. How should the response profiles of these more "heterogeneous" cells be described, and how do they contribute to behavior? In this work, we took a computational approach to addressing these questions. We first performed a statistical analysis that shows that heterogeneous MEC cells are just as reliable in their response patterns as the more stereotypical cell types, suggesting that they have a coherent functional role. Next, we evaluated a spectrum of candidate models in terms of their ability to describe the response profiles of both stereotypical and heterogeneous MEC cells. We found that recently developed task-optimized neural network models are substantially better than traditional grid cell-centric models at matching most MEC neuronal response profiles --- including those of grid cells themselves --- despite not being explicitly trained for this purpose. Specific choices of network architecture (such as gated nonlinearities and an explicit intermediate place cell representation) have an important effect on the ability of the model to generalize to novel scenarios, with the best of these models closely approaching the noise ceiling of the data itself. We then performed in silico experiments on this model to address questions involving the relative functional relevance of various cell types, finding that heterogeneous cells are likely to be just as involved in downstream functional outcomes (such as path integration) as grid and border cells. Finally, inspired by recent data showing that, going beyond their spatial response selectivity, MEC cells are also responsive to non-spatial rewards, we introduce a new MEC model that performs reward-modulated path integration. We find that this unified model matches neural recordings across all variable-reward conditions. Taken together, our results point toward a conceptually principled goal-driven modeling approach for moving future experimental and computational efforts beyond overly-simplistic single-cell stereotypes.

[Beyond Smoothness: Incorporating Low-Rank Analysis into Nonparametric Density Estimation](#)

- Robert A. Vandermeulen, Antoine Ledent
- abstract@[open-review](#): The construction and theoretical analysis of the most popular universally consistent nonparametric density estimators hinge on one functional property: smoothness. In this paper we investigate the theoretical implications of incorporating a multi-view latent variable model, a type of low-rank model, into nonparametric density estimation. To do this we perform extensive analysis on histogram-style estimators that integrate a multi-view model. Our analysis culminates in showing that there exists a universally consistent histogram-style estimator that converges to any multi-view model with a finite number of Lipschitz continuous components at a rate of $\widetilde{O}(1/\sqrt{3n})$ in L^1 error. In contrast, the standard histogram estimator can converge at a rate slower than $1/\sqrt{d}n$ on the same class of densities. We also introduce a new nonparametric latent variable model based on the Tucker decomposition. A rudimentary implementation of our estimators experimentally demonstrates a considerable performance improvement over the standard histogram estimator. We also provide a thorough analysis of the sample complexity of our Tucker decomposition-based model and a variety of other results. Thus, our paper provides solid theoretical foundations for extending low-rank techniques to the nonparametric setting.

[Multi-View Representation Learning via Total Correlation Objective](#)

- HyeongJoo Hwang, Geon-Hyeong Kim, Seunghoon Hong, Kee-Eung Kim
- abstract@[open-review](#): Multi-View Representation Learning (MVRL) aims to discover a shared representation of observations from different views with the complex underlying correlation. In this paper, we propose a variational approach which casts MVRL as maximizing the amount of total correlation reduced by the representation, aiming to learn a shared latent representation that is informative yet succinct to capture the correlation among multiple views. To this end, we introduce a tractable surrogate objective function under the proposed framework, which allows our method to fuse and calibrate the observations in the representation space. From the information-theoretic perspective, we show that our framework subsumes existing multi-view generative models. Lastly, we show that our approach straightforwardly extends to the Partial MVRL (PMVRL) setting, where the observations are missing without any regular pattern. We demonstrate the effectiveness of our approach in the multi-view translation and classification tasks, outperforming strong baseline methods.

[FACMAC: Factored Multi-Agent Centralised Policy Gradients](#)

- Bei Peng, Tabish Rashid, Christian Schroeder de Witt, Pierre-Alexandre Kamienny, Philip Torr, Wendelin Boehmer, Shimon Whiteson
- abstract@[open-review](#): We propose FACtored Multi-Agent Centralised policy gradients (FACMAC), a new method for cooperative multi-agent reinforcement learning in both discrete and continuous action spaces. Like MADDPG, a popular multi-agent actor-critic method, our approach uses deep

deterministic policy gradients to learn policies. However, FACMAC learns a centralised but factored critic, which combines per-agent utilities into the joint action-value function via a non-linear monotonic function, as in QMIX, a popular multi-agent $\$Q\$$ -learning algorithm. However, unlike QMIX, there are no inherent constraints on factoring the critic. We thus also employ a nonmonotonic factorisation and empirically demonstrate that its increased representational capacity allows it to solve some tasks that cannot be solved with monolithic, or monotonically factored critics. In addition, FACMAC uses a centralised policy gradient estimator that optimises over the entire joint action space, rather than optimising over each agent's action space separately as in MADDPG. This allows for more coordinated policy changes and fully reaps the benefits of a centralised critic. We evaluate FACMAC on variants of the multi-agent particle environments, a novel multi-agent MuJoCo benchmark, and a challenging set of StarCraft II micromanagement tasks. Empirical results demonstrate FACMAC's superior performance over MADDPG and other baselines on all three domains.

[EDGE: Explaining Deep Reinforcement Learning Policies](#)

- Wenbo Guo, Xian Wu, Usmann Khan, Xinyu Xing
- abstract@[open-review](#): With the rapid development of deep reinforcement learning (DRL) techniques, there is an increasing need to understand and interpret DRL policies. While recent research has developed explanation methods to interpret how an agent determines its moves, they cannot capture the importance of actions/states to a game's final result. In this work, we propose a novel self-explainable model that augments a Gaussian process with a customized kernel function and an interpretable predictor. Together with the proposed model, we also develop a parameter learning procedure that leverages inducing points and variational inference to improve learning efficiency. Using our proposed model, we can predict an agent's final rewards from its game episodes and extract time step importance within episodes as strategy-level explanations for that agent. Through experiments on Atari and MuJoCo games, we verify the explanation fidelity of our method and demonstrate how to employ interpretation to understand agent behavior, discover policy vulnerabilities, remediate policy errors, and even defend against adversarial attacks.

[Learning to Assimilate in Chaotic Dynamical Systems](#)

- John McCabe, Jed Brown
- abstract@[open-review](#): The accuracy of simulation-based forecasting in chaotic systems is heavily dependent on high-quality estimates of the system state at the beginning of the forecast. Data assimilation methods are used to infer these initial conditions by systematically combining noisy, incomplete observations and numerical models of system dynamics to produce highly effective estimation schemes. We introduce a self-supervised framework, which we call \textit{amortized assimilation}, for learning to assimilate in dynamical systems. Amortized assimilation combines deep learning-based denoising with differentiable simulation, using independent neural networks to assimilate specific observation types while connecting the gradient flow between these sub-tasks with differentiable simulation and shared recurrent memory. This hybrid architecture admits a self-supervised training objective which is minimized by an unbiased estimator of the true system state even in the presence of only noisy training data. Numerical experiments across several chaotic benchmark systems highlight the improved effectiveness of our approach compared to widely-used data assimilation methods.

[Object-aware Contrastive Learning for Debiased Scene Representation](#)

- Sangwoo Mo, Hyunwoo Kang, Kihyuk Sohn, Chun-Liang Li, Jinwoo Shin
- abstract@[open-review](#): Contrastive self-supervised learning has shown impressive results in learning visual representations from unlabeled images by enforcing invariance against different data augmentations. However, the learned representations are often contextually biased to the spurious scene correlations of different objects or object and background, which may harm their generalization on the downstream tasks. To tackle the issue, we develop a novel object-aware contrastive learning framework that first (a) localizes objects in a self-supervised manner and then (b) debias scene correlations via appropriate data augmentations considering the inferred object locations. For (a), we propose the contrastive class activation map (ContraCAM), which finds the most discriminative regions (e.g., objects) in the image compared to the other images using the contrastively trained models. We further improve the ContraCAM to detect multiple objects and entire shapes via an iterative refinement procedure. For (b), we introduce two data augmentations based on ContraCAM, object-aware random crop and background mixup, which reduce contextual and background biases during contrastive self-supervised learning, respectively. Our experiments demonstrate the effectiveness of our representation learning framework, particularly when trained under multi-object images or evaluated under the background (and distribution) shifted images. Code is available at <https://github.com/alinlab/object-aware-contrastive>.

[Evaluating Efficient Performance Estimators of Neural Architectures](#)

- Xuefei Ning, Changcheng Tang, Wenshuo Li, Zixuan Zhou, Shuang Liang, Huazhong Yang, Yu Wang
- abstract@[open-review](#): Conducting efficient performance estimations of neural architectures is a major challenge in neural architecture search (NAS). To reduce the architecture training costs in NAS, one-shot estimators (OSEs) amortize the architecture training costs by sharing the parameters of one supernet between all architectures. Recently, zero-shot estimators (ZSEs) that involve no training are proposed to further reduce the architecture evaluation cost. Despite the high efficiency of these estimators, the quality of such estimations has not been thoroughly studied. In this paper, we conduct an extensive and organized assessment of OSEs and ZSEs on five NAS benchmarks: NAS-Bench-101/201/301, and NDS ResNet/ResNeXt-A. Specifically, we employ a set of NAS-oriented criteria to study the behavior of OSEs and ZSEs, and reveal their biases and variances. After analyzing how and why the OSE estimations are unsatisfying, we explore how to mitigate the correlation gap of OSEs from three perspectives. Through our analysis, we give out suggestions for future application and development of efficient architecture performance estimators. Furthermore, the analysis framework proposed in our work could be utilized in future research to give a more comprehensive understanding of newly designed architecture performance estimators. The code is available at https://github.com/walkerning/aw_nas.

[A-NeRF: Articulated Neural Radiance Fields for Learning Human Shape, Appearance, and Pose](#)

- Shih-Yang Su, Frank Yu, Michael Zollhoefer, Helge Rhodin
- abstract@[open-review](#): While deep learning reshaped the classical motion capture pipeline with feed-forward networks, generative models are required to recover fine alignment via iterative refinement. Unfortunately, the existing models are usually hand-crafted or learned in controlled conditions, only applicable to limited domains. We propose a method to learn a generative neural body model from unlabelled monocular videos by extending Neural Radiance Fields (NeRFs). We equip them with a skeleton to apply to time-varying and articulated motion. A key insight is that implicit models require the inverse of the forward kinematics used in explicit surface models. Our reparameterization defines spatial latent variables relative to the pose of body parts and thereby overcomes ill-posed inverse operations with an overparameterization. This enables learning volumetric body shape and appearance from scratch while jointly refining the articulated pose; all without ground truth labels for appearance, pose, or 3D shape on the input videos. When used for novel-view-synthesis and motion capture, our neural model improves accuracy on diverse datasets.

[Differential Privacy Over Riemannian Manifolds](#)

- Matthew Reimherr, Karthik Bharath, Carlos Soto
- abstract@[open-review](#): In this work we consider the problem of releasing a differentially private statistical summary that resides on a Riemannian manifold. We present an extension of the Laplace or K-norm mechanism that utilizes intrinsic distances and volumes on the manifold. We also consider in detail the specific case where the summary is the Fréchet mean of data residing on a manifold. We demonstrate that our mechanism is rate optimal and depends only on the dimension of the manifold, not on the dimension of any ambient space, while also showing how ignoring the manifold structure can

decrease the utility of the sanitized summary. We illustrate our framework in two examples of particular interest in statistics: the space of symmetric positive definite matrices, which is used for covariance matrices, and the sphere, which can be used as a space for modeling discrete distributions.

[How can classical multidimensional scaling go wrong?](#)

- Rishi Sonthalia, Greg Van Buskirk, Benjamin Raichel, Anna Gilbert
- abstract@[open-review](#): Given a matrix D describing the pairwise dissimilarities of a data set, a common task is to embed the data points into Euclidean space. The classical multidimensional scaling (cMDS) algorithm is a widespread method to do this. However, theoretical analysis of the robustness of the algorithm and an in-depth analysis of its performance on non-Euclidean metrics is lacking. In this paper, we derive a formula, based on the eigenvalues of a matrix obtained from D , for the Frobenius norm of the difference between D and the metric D_{cmds} returned by cMDS. This error analysis leads us to the conclusion that when the derived matrix has a significant number of negative eigenvalues, then $\|D - D_{\text{cmds}}\|_F$, after initially decreasing, will eventually increase as we increase the dimension. Hence, counterintuitively, the quality of the embedding degrades as we increase the dimension. We empirically verify that the Frobenius norm increases as we increase the dimension for a variety of non-Euclidean metrics. We also show on several benchmark datasets that this degradation in the embedding results in the classification accuracy of both simple (e.g., 1-nearest neighbor) and complex (e.g., multi-layer neural nets) classifiers decreasing as we increase the embedding dimension. Finally, our analysis leads us to a new efficiently computable algorithm that returns a matrix D_l that is at least as close to the original distances as D_t (the Euclidean metric closest in ℓ_2 distance). While D_l is not metric, when given as input to cMDS instead of D , it empirically results in solutions whose distance to D does not increase when we increase the dimension and the classification accuracy degrades less than the cMDS solution.

[Modeling Heterogeneous Hierarchies with Relation-specific Hyperbolic Cones](#)

- Yushi Bai, Zhitao Ying, Hongyu Ren, Jure Leskovec
- abstract@[open-review](#): Hierarchical relations are prevalent and indispensable for organizing human knowledge captured by a knowledge graph (KG). The key property of hierarchical relations is that they induce a partial ordering over the entities, which needs to be modeled in order to allow for hierarchical reasoning. However, current KG embeddings can model only a single global hierarchy (single global partial ordering) and fail to model multiple heterogeneous hierarchies that exist in a single KG. Here we present ConE (Cone Embedding), a KG embedding model that is able to simultaneously model multiple hierarchical as well as non-hierarchical relations in a knowledge graph. ConE embeds entities into hyperbolic cones and models relations as transformations between the cones. In particular, ConE uses cone containment constraints in different subspaces of the hyperbolic embedding space to capture multiple heterogeneous hierarchies. Experiments on standard knowledge graph benchmarks show that ConE obtains state-of-the-art performance on hierarchical reasoning tasks as well as knowledge graph completion task on hierarchical graphs. In particular, our approach yields new state-of-the-art Hits@1 of 45.3% on WN18RR and 16.1% on DDB14 (0.231 MRR). As for hierarchical reasoning task, our approach outperforms previous best results by an average of 20% across the three datasets.

[Non-asymptotic Error Bounds for Bidirectional GANs](#)

- Shiao Liu, Yunfei Yang, Jian Huang, Yuling Jiao, Yang Wang
- abstract@[open-review](#): We derive nearly sharp bounds for the bidirectional GAN (BiGAN) estimation error under the Dudley distance between the latent joint distribution and the data joint distribution with appropriately specified architecture of the neural networks used in the model. To the best of our knowledge, this is the first theoretical guarantee for the bidirectional GAN learning approach. An appealing feature of our results is that they do not assume the reference and the data distributions to have the same dimensions or these distributions to have bounded support. These assumptions are commonly assumed in the existing convergence analysis of the unidirectional GANs but may not be satisfied in practice. Our results are also applicable to the Wasserstein bidirectional GAN if the target distribution is assumed to have a bounded support. To prove these results, we construct neural network functions that push forward an empirical distribution to another arbitrary empirical distribution on a possibly different-dimensional space. We also develop a novel decomposition of the integral probability metric for the error analysis of bidirectional GANs. These basic theoretical results are of independent interest and can be applied to other related learning problems.

[Confidence-Aware Imitation Learning from Demonstrations with Varying Optimality](#)

- Songyuan Zhang, ZHANGJIE CAO, Dorsa Sadigh, Yanan Sui
- abstract@[open-review](#): Most existing imitation learning approaches assume the demonstrations are drawn from experts who are optimal, but relaxing this assumption enables us to use a wider range of data. Standard imitation learning may learn a suboptimal policy from demonstrations with varying optimality. Prior works use confidence scores or rankings to capture beneficial information from demonstrations with varying optimality, but they suffer from many limitations, e.g., manually annotated confidence scores or high average optimality of demonstrations. In this paper, we propose a general framework to learn from demonstrations with varying optimality that jointly learns the confidence score and a well-performing policy. Our approach, Confidence-Aware Imitation Learning (CAIL) learns a well-performing policy from confidence-reweighted demonstrations, while using an outer loss to track the performance of our model and to learn the confidence. We provide theoretical guarantees on the convergence of CAIL and evaluate its performance in both simulated and real robot experiments. Our results show that CAIL significantly outperforms other imitation learning methods from demonstrations with varying optimality. We further show that even without access to any optimal demonstrations, CAIL can still learn a successful policy, and outperforms prior work.

[Answering Complex Causal Queries With the Maximum Causal Set Effect](#)

- Zachary Markovich
- abstract@[open-review](#): The standard tools of causal inference have been developed to answer simple causal queries which can be easily formalized as a small number of statistical estimands in the context of a particular structural causal model (SCM); however, scientific theories often make diffuse predictions about a large number of causal variables. This article proposes a framework for parameterizing such complex causal queries as the maximum difference in causal effects associated with two sets of causal variables that have a researcher specified probability of occurring. We term this estimand the Maximum Causal Set Effect (MCSE) and develop an estimator for it that is asymptotically consistent and conservative in finite samples under assumptions that are standard in the causal inference literature. This estimator is also asymptotically normal and amenable to the non-parametric bootstrap, facilitating classical statistical inference about this novel estimand. We compare this estimator to more common latent variable approaches and find that it can uncover larger causal effects in both real world and simulated data.

[Identifiability in inverse reinforcement learning](#)

- Haoyang Cao, Samuel Cohen, Lukasz Szpruch
- abstract@[open-review](#): Inverse reinforcement learning attempts to reconstruct the reward function in a Markov decision problem, using observations of agent actions. As already observed in Russell [1998] the problem is ill-posed, and the reward function is not identifiable, even under the presence of perfect information about optimal behavior. We provide a resolution to this non-identifiability for problems with entropy regularization. For a given environment, we fully characterize the reward functions leading to a given policy and demonstrate that, given demonstrations of actions for the same reward under two distinct discount factors, or under sufficiently different environments, the unobserved reward can be recovered up to a constant. We also give general necessary and sufficient conditions for reconstruction of time-homogeneous rewards on finite horizons, and for action-independent rewards, generalizing recent results of Kim et al. [2021] and Fu et al. [2018].

[A Probabilistic State Space Model for Joint Inference from Differential Equations and Data](#)

- Jonathan Schmidt, Nicholas Krämer, Philipp Hennig
- abstract@[open-review](#): Mechanistic models with differential equations are a key component of scientific applications of machine learning. Inference in such models is usually computationally demanding because it involves repeatedly solving the differential equation. The main problem here is that the numerical solver is hard to combine with standard inference techniques. Recent work in probabilistic numerics has developed a new class of solvers for ordinary differential equations (ODEs) that phrase the solution process directly in terms of Bayesian filtering. We here show that this allows such methods to be combined very directly, with conceptual and numerical ease, with latent force models in the ODE itself. It then becomes possible to perform approximate Bayesian inference on the latent force as well as the ODE solution in a single, linear complexity pass of an extended Kalman filter / smoother — that is, at the cost of computing a single ODE solution. We demonstrate the expressiveness and performance of the algorithm by training, among others, a non-parametric SIRD model on data from the COVID-19 outbreak.

[On Plasticity, Invariance, and Mutually Frozen Weights in Sequential Task Learning](#)

- Julian Zilly, Alessandro Achille, Andrea Censi, Emilio Frazzoli
- abstract@[open-review](#): Plastic neural networks have the ability to adapt to new tasks. However, in a continual learning setting, the configuration of parameters learned in previous tasks can severely reduce the adaptability to future tasks. In particular, we show that, when using weight decay, weights in successive layers of a deep network may become "mutually frozen". This has a double effect: on the one hand, it makes the network updates more invariant to nuisance factors, providing a useful bias for future tasks. On the other hand, it can prevent the network from learning new tasks that require significantly different features. In this context, we find that the local input sensitivity of a deep model is correlated with its ability to adapt, thus leading to an intriguing trade-off between adaptability and invariance when training a deep model more than once. We then show that a simple intervention that "resets" the mutually frozen connections can improve transfer learning on a variety of visual classification tasks. The efficacy of "resetting" itself depends on the size of the target dataset and the difference of the pre-training and target domains, allowing us to achieve state-of-the-art results on some datasets.

[Provably Efficient Black-Box Action Poisoning Attacks Against Reinforcement Learning](#)

- Guanlin Liu, Lifeng LAI
- abstract@[open-review](#): Due to the broad range of applications of reinforcement learning (RL), understanding the effects of adversarial attacks against RL model is essential for the safe applications of this model. Prior theoretical works on adversarial attacks against RL mainly focus on either reward poisoning attacks or environment poisoning attacks. In this paper, we introduce a new class of attacks named action poisoning attacks, where an adversary can change the action signal selected by the agent. Compared with existing attack models, the attacker's ability in the proposed action poisoning attack model is more restricted, which brings some design challenges. We study the action poisoning attack in both white-box and black-box settings. We introduce an adaptive attack scheme called LCB-H, which works for most RL agents in the black-box setting. We prove that LCB-H attack can force any efficient RL agent, whose dynamic regret scales sublinearly with the total number of steps taken, to choose actions according to a policy selected by the attacker very frequently, with only sublinear cost. In addition, we apply LCB-H attack against a very popular model-free RL algorithm: UCB-H. We show that, even in black-box setting, by spending only logarithm cost, the proposed LCB-H attack scheme can force the UCB-H agent to choose actions according to the policy selected by the attacker very frequently.

[Fast Approximation of the Sliced-Wasserstein Distance Using Concentration of Random Projections](#)

- Kimia Nadjahi, Alain Durmus, Pierre E Jacob, Roland Badeau, Umut Simsekli
- abstract@[open-review](#): The Sliced-Wasserstein distance (SW) is being increasingly used in machine learning applications as an alternative to the Wasserstein distance and offers significant computational and statistical benefits. Since it is defined as an expectation over random projections, SW is commonly approximated by Monte Carlo. We adopt a new perspective to approximate SW by making use of the concentration of measure phenomenon: under mild assumptions, one-dimensional projections of a high-dimensional random vector are approximately Gaussian. Based on this observation, we develop a simple deterministic approximation for SW. Our method does not require sampling a number of random projections, and is therefore both accurate and easy to use compared to the usual Monte Carlo approximation. We derive nonasymptotical guarantees for our approach, and show that the approximation error goes to zero as the dimension increases, under a weak dependence condition on the data distribution. We validate our theoretical findings on synthetic datasets, and illustrate the proposed approximation on a generative modeling problem.

[Causal Navigation by Continuous-time Neural Networks](#)

- Charles Vorbach, Ramin Hasani, Alexander Amini, Mathias Lechner, Daniela Rus
- abstract@[open-review](#): Imitation learning enables high-fidelity, vision-based learning of policies within rich, photorealistic environments. However, such techniques often rely on traditional discrete-time neural models and face difficulties in generalizing to domain shifts by failing to account for the causal relationships between the agent and the environment. In this paper, we propose a theoretical and experimental framework for learning causal representations using continuous-time neural networks, specifically over their discrete-time counterparts. We evaluate our method in the context of visual-control learning of drones over a series of complex tasks, ranging from short- and long-term navigation, to chasing static and dynamic objects through photorealistic environments. Our results demonstrate that causal continuous-time deep models can perform robust navigation tasks, where advanced recurrent models fail. These models learn complex causal control representations directly from raw visual inputs and scale to solve a variety of tasks using imitation learning.

[Global Convergence of Online Optimization for Nonlinear Model Predictive Control](#)

- Sen Na
- abstract@[open-review](#): We study a real-time iteration (RTI) scheme for solving online optimization problem appeared in nonlinear optimal control. The proposed RTI scheme modifies the existing RTI-based model predictive control (MPC) algorithm, by selecting the stepsize of each Newton step at each sampling time using a differentiable exact augmented Lagrangian. The scheme can adaptively select the penalty parameters of augmented Lagrangian on the fly, which are shown to be stabilized after certain time periods. We prove under generic assumptions that, by involving stepsize selection instead of always using a full Newton step (like what most of the existing RTIs do), the scheme converges globally: for any initial point, the KKT residuals of the subproblems converge to zero. A key step is to show that augmented Lagrangian keeps decreasing as horizon moves forward. We demonstrate the global convergence behavior of the proposed RTI scheme in a numerical experiment.

[Argmax Flows and Multinomial Diffusion: Learning Categorical Distributions](#)

- Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, Max Welling
- abstract@[open-review](#): Generative flows and diffusion models have been predominantly trained on ordinal data, for example natural images. This paper introduces two extensions of flows and diffusion for categorical data such as language or image segmentation: Argmax Flows and Multinomial Diffusion. Argmax Flows are defined by a composition of a continuous distribution (such as a normalizing flow), and an argmax function. To optimize this model, we learn a probabilistic inverse for the argmax that lifts the categorical data to a continuous space. Multinomial Diffusion gradually adds categorical noise

in a diffusion process, for which the generative denoising process is learned. We demonstrate that our method outperforms existing dequantization approaches on text modelling and modelling on image segmentation maps in log-likelihood.

[Learning with User-Level Privacy](#)

- Daniel Levy, Ziteng Sun, Kareem Amin, Satyen Kale, Alex Kulesza, Mehryar Mohri, Ananda Theertha Suresh
- abstract@[open-review](#): We propose and analyze algorithms to solve a range of learning tasks under user-level differential privacy constraints. Rather than guaranteeing only the privacy of individual samples, user-level DP protects a user's entire contribution ($\$m \geq 1\$$ samples), providing more stringent but more realistic protection against information leaks. We show that for high-dimensional mean estimation, empirical risk minimization with smooth losses, stochastic convex optimization, and learning hypothesis classes with finite metric entropy, the privacy cost decreases as $\$O(1/\sqrt{m})\$$ as users provide more samples. In contrast, when increasing the number of users $\$n\$$, the privacy cost decreases at a faster $\$O(1/n)\$$ rate. We complement these results with lower bounds showing the minimax optimality of our algorithms for mean estimation and stochastic convex optimization. Our algorithms rely on novel techniques for private mean estimation in arbitrary dimension with error scaling as the concentration radius $\$|\tau|$ of the distribution rather than the entire range.

[Don't Generate Me: Training Differentially Private Generative Models with Sinkhorn Divergence](#)

- Tianshi Cao, Alex Bie, Arash Vahdat, Sanja Fidler, Karsten Kreis
- abstract@[open-review](#): Although machine learning models trained on massive data have led to breakthroughs in several areas, their deployment in privacy-sensitive domains remains limited due to restricted access to data. Generative models trained with privacy constraints on private data can sidestep this challenge, providing indirect access to private data instead. We propose DP-Sinkhorn, a novel optimal transport-based generative method for learning data distributions from private data with differential privacy. DP-Sinkhorn minimizes the Sinkhorn divergence, a computationally efficient approximation to the exact optimal transport distance, between the model and data in a differentially private manner and uses a novel technique for controlling the bias-variance trade-off of gradient estimates. Unlike existing approaches for training differentially private generative models, which are mostly based on generative adversarial networks, we do not rely on adversarial objectives, which are notoriously difficult to optimize, especially in the presence of noise imposed by privacy constraints. Hence, DP-Sinkhorn is easy to train and deploy. Experimentally, we improve upon the state-of-the-art on multiple image modeling benchmarks and show differentially private synthesis of informative RGB images.

[Keeping Your Eye on the Ball: Trajectory Attention in Video Transformers](#)

- Mandela Patrick, Dylan Campbell, Yuki Asano, Ishan Misra, Florian Metze, Christoph Feichtenhofer, Andrea Vedaldi, João F. Henriques
- abstract@[open-review](#): In video transformers, the time dimension is often treated in the same way as the two spatial dimensions. However, in a scene where objects or the camera may move, a physical point imaged at one location in frame $\$t\$$ may be entirely unrelated to what is found at that location in frame $\$t+k\$$. These temporal correspondences should be modeled to facilitate learning about dynamic scenes. To this end, we propose a new drop-in block for video transformers - trajectory attention - that aggregates information along implicitly determined motion paths. We additionally propose a new method to address the quadratic dependence of computation and memory on the input size, which is particularly important for high resolution or long videos. While these ideas are useful in a range of settings, we apply them to the specific task of video action recognition with a transformer model and obtain state-of-the-art results on the Kinetics, Something-Something V2, and Epic-Kitchens datasets.

[Variational Bayesian Optimistic Sampling](#)

- Brendan O'Donoghue, Tor Lattimore
- abstract@[open-review](#): We consider online sequential decision problems where an agent must balance exploration and exploitation. We derive a set of 'Bayesian optimistic' policies which, in the stochastic multi-armed bandit case, includes the Thompson sampling policy. We provide a new analysis showing that any algorithm producing policies in the optimistic set enjoys $\$|\tilde{O}(\sqrt{AT})|\$$ Bayesian regret for a problem with $\$A\$$ actions after $\$T\$$ rounds. We extend the regret analysis for optimistic policies to bilinear saddle-point problems which include zero-sum matrix games and constrained bandits as special cases. In this case we show that Thompson sampling can produce policies outside of the optimistic set and suffer linear regret in some instances. Finding a policy inside the optimistic set amounts to solving a convex optimization problem and we call the resulting algorithm 'variational Bayesian optimistic sampling' (VBOS). The procedure works for any posteriors, i.e., it does not require the posterior to have any special properties, such as log-concavity, unimodality, or smoothness. The variational view of the problem has many useful properties, including the ability to tune the exploration-exploitation tradeoff, add regularization, incorporate constraints, and linearly parameterize the policy.

[Cross-modal Domain Adaptation for Cost-Efficient Visual Reinforcement Learning](#)

- Xiong-Hui Chen, Shengyi Jiang, Feng Xu, Zongzhang Zhang, Yang Yu
- abstract@[open-review](#): In visual-input sim-to-real scenarios, to overcome the reality gap between images rendered in simulators and those from the real world, domain adaptation, i.e., learning an aligned representation space between simulators and the real world, then training and deploying policies in the aligned representation, is a promising direction. Previous methods focus on same-modal domain adaptation. However, those methods require building and running simulators that render high-quality images, which can be difficult and costly. In this paper, we consider a more cost-efficient setting of visual-input sim-to-real where only low-dimensional states are simulated. We first point out that the objective of learning mapping functions in previous methods that align the representation spaces is ill-posed, prone to yield an incorrect mapping. When the mapping crosses modalities, previous methods are easier to fail. Our algorithm, Cross-mOdal Domain Adaptation with Sequential structure (CODAS), mitigates the ill-posedness by utilizing the sequential nature of the data sampling process in RL tasks. Experiments on MuJoCo and Hand Manipulation Suite tasks show that the agents deployed with our method achieve similar performance as it has in the source domain, while those deployed with previous methods designed for same-modal domain adaptation suffer a larger performance gap.

[D2C: Diffusion-Decoding Models for Few-Shot Conditional Generation](#)

- Abhishek Sinha, Jiaming Song, Chenlin Meng, Stefano Ermon
- abstract@[open-review](#): Conditional generative models of high-dimensional images have many applications, but supervision signals from conditions to images can be expensive to acquire. This paper describes Diffusion-Decoding models with Contrastive representations (D2C), a paradigm for training unconditional variational autoencoders (VAE) for few-shot conditional image generation. D2C uses a learned diffusion-based prior over the latent representations to improve generation and contrastive self-supervised learning to improve representation quality. D2C can adapt to novel generation tasks, conditioned on labels or manipulation constraints, by learning from as few as 100 labeled examples. On conditional generation from new labels, D2C achieves superior performance over state-of-the-art VAEs and diffusion models. On conditional image manipulation, D2C generations are two orders of magnitude faster to produce over StyleGAN2 ones and are preferred by 50% - 60% of the human evaluators in a double-blind study. We release our code at <https://github.com/jiamings/d2c>.

[Continual Auxiliary Task Learning](#)

- Matthew McLeod, Chunlok Lo, Matthew Schlegel, Andrew Jacobsen, Raksha Kumaraswamy, Martha White, Adam White
- abstract@[open-review](#): Learning auxiliary tasks, such as multiple predictions about the world, can provide many benefits to reinforcement learning systems. A variety of off-policy learning algorithms have been developed to learn such predictions, but as yet there is little work on how to adapt the behavior to gather useful data for those off-policy predictions. In this work, we investigate a reinforcement learning system designed to learn a collection of auxiliary tasks, with a behavior policy learning to take actions to improve those auxiliary predictions. We highlight the inherent non-stationarity in this continual auxiliary task learning problem, for both prediction learners and the behavior learner. We develop an algorithm based on successor features that facilitates tracking under non-stationary rewards, and prove the separation into learning successor features and rewards provides convergence rate improvements. We conduct an in-depth study into the resulting multi-prediction learning system.

[Constrained Two-step Look-Ahead Bayesian Optimization](#)

- Yunxiang Zhang, Xiangyu Zhang, Peter Frazier
- abstract@[open-review](#): Recent advances in computationally efficient non-myopic Bayesian optimization offer improved query efficiency over traditional myopic methods like expected improvement, with only a modest increase in computational cost. These advances have been largely limited to unconstrained BO methods with only a few exceptions which require heavy computation. For instance, one existing multi-step lookahead constrained BO method (Lam & Willcox, 2017) relies on computationally expensive unreliable brute force derivative-free optimization of a Monte Carlo rollout acquisition function. Methods that use the reparameterization trick for more efficient derivative-based optimization of non-myopic acquisition functions in the unconstrained setting, like sample average approximation and infinitesimal perturbation analysis, do not extend: constraints introduce discontinuities in the sampled acquisition function surface. Moreover, we argue here that being non-myopic is even more important in constrained problems because fear of violating constraints pushes myopic methods away from sampling the boundary between feasible and infeasible regions, slowing the discovery of optimal solutions with tight constraints. In this paper, we propose a computationally efficient two-step lookahead constrained Bayesian optimization acquisition function (2-OPT-C) supporting both sequential and batch settings. To enable fast acquisition function optimization, we develop a novel likelihood ratio-based unbiased estimator of the gradient of the two-step optimal acquisition function that does not use the reparameterization trick. In numerical experiments, 2-OPT-C typically improves query efficiency by 2x or more over previous methods, and in some cases by 10x or more.

[Learning with Labeling Induced Abstentions](#)

- Kareem Amin, Giulia DeSalvo, Afshin Rostamizadeh
- abstract@[open-review](#): Consider a setting where we wish to automate an expensive task with a machine learning algorithm using a limited labeling resource. In such settings, examples routed for labeling are often out of scope for the machine learning algorithm. For example, in a spam detection setting, human reviewers not only provide labeled data but are such high-quality detectors of spam that examples routed to them no longer require machine evaluation. As a consequence, the distribution of examples routed to the machine is intimately tied to the process generating labels. We introduce a formalization of this setting, and give an algorithm that simultaneously learns a model and decides when to request a label by leveraging ideas from both the abstention and active learning literatures. We prove an upper bound on the algorithm's label complexity and a matching lower bound for any algorithm in this setting. We conduct a thorough set of experiments including an ablation study to test different components of our algorithm. We demonstrate the effectiveness of an efficient version of our algorithm over margin sampling on a variety of datasets.

[SQALER: Scaling Question Answering by Decoupling Multi-Hop and Logical Reasoning](#)

- Mattia Atzeni, Jasmina Bogojeska, Andreas Loukas
- abstract@[open-review](#): State-of-the-art approaches to reasoning and question answering over knowledge graphs (KGs) usually scale with the number of edges and can only be applied effectively on small instance-dependent subgraphs. In this paper, we address this issue by showing that multi-hop and more complex logical reasoning can be accomplished separately without losing expressive power. Motivated by this insight, we propose an approach to multi-hop reasoning that scales linearly with the number of relation types in the graph, which is usually significantly smaller than the number of edges or nodes. This produces a set of candidate solutions that can be provably refined to recover the solution to the original problem. Our experiments on knowledge-based question answering show that our approach solves the multi-hop MetaQA dataset, achieves a new state-of-the-art on the more challenging WebQuestionsSP, is orders of magnitude more scalable than competitive approaches, and can achieve compositional generalization out of the training distribution.

[Out-of-Distribution Generalization in Kernel Regression](#)

- Abdulkadir Canatar, Blake Bordelon, Cengiz Pehlevan
- abstract@[open-review](#): In real word applications, data generating process for training a machine learning model often differs from what the model encounters in the test stage. Understanding how and whether machine learning models generalize under such distributional shifts have been a theoretical challenge. Here, we study generalization in kernel regression when the training and test distributions are different using methods from statistical physics. Using the replica method, we derive an analytical formula for the out-of-distribution generalization error applicable to any kernel and real datasets. We identify an overlap matrix that quantifies the mismatch between distributions for a given kernel as a key determinant of generalization performance under distribution shift. Using our analytical expressions we elucidate various generalization phenomena including possible improvement in generalization when there is a mismatch. We develop procedures for optimizing training and test distributions for a given data budget to find best and worst case generalizations under the shift. We present applications of our theory to real and synthetic datasets and for many kernels. We compare results of our theory applied to Neural Tangent Kernel with simulations of wide networks and show agreement. We analyze linear regression in further depth.

[FL-WBC: Enhancing Robustness against Model Poisoning Attacks in Federated Learning from a Client Perspective](#)

- Jingwei Sun, Ang Li, Louis DiValentin, Amin Hassanzadeh, Yiran Chen, Hai Li
- abstract@[open-review](#): Federated learning (FL) is a popular distributed learning framework that trains a global model through iterative communications between a central server and edge devices. Recent works have demonstrated that FL is vulnerable to model poisoning attacks. Several server-based defense approaches (e.g. robust aggregation), have been proposed to mitigate such attacks. However, we empirically show that under extremely strong attacks, these defensive methods fail to guarantee the robustness of FL. More importantly, we observe that as long as the global model is polluted, the impact of attacks on the global model will remain in subsequent rounds even if there are no subsequent attacks. In this work, we propose a client-based defense, named White Blood Cell for Federated Learning (FL-WBC), which can mitigate model poisoning attacks that have already polluted the global model. The key idea of FL-WBC is to identify the parameter space where long-lasting attack effect on parameters resides and perturb that space during local training. Furthermore, we derive a certified robustness guarantee against model poisoning attacks and a convergence guarantee to FedAvg after applying our FL-WBC. We conduct experiments on FashionMNIST and CIFAR10 to evaluate the defense against state-of-the-art model poisoning attacks. The results demonstrate that our method can effectively mitigate model poisoning attack impact on the global model within 5 communication rounds with nearly no accuracy drop under both IID and Non-IID settings. Our defense is also complementary to existing server-based robust aggregation approaches and can further improve the robustness of FL under extremely strong attacks.

[Chebyshev-Cantelli PAC-Bayes-Bennett Inequality for the Weighted Majority Vote](#)

- Yi-Shan Wu, Andres Masegosa, Stephan Lorenzen, Christian Igel, Yevgeny Seldin

- abstract@[open-review](#): We present a new second-order oracle bound for the expected risk of a weighted majority vote. The bound is based on a novel parametric form of the Chebyshev-Cantelli inequality (a.k.a. one-sided Chebyshev's), which is amenable to efficient minimization. The new form resolves the optimization challenge faced by prior oracle bounds based on the Chebyshev-Cantelli inequality, the C-bounds [Germain et al., 2015], and, at the same time, it improves on the oracle bound based on second order Markov's inequality introduced by Masegosa et al. [2020]. We also derive a new concentration of measure inequality, which we name PAC-Bayes-Bennett, since it combines PAC-Bayesian bounding with Bennett's inequality. We use it for empirical estimation of the oracle bound. The PAC-Bayes-Bennett inequality improves on the PAC-Bayes-Bernstein inequality of Seldin et al. [2012]. We provide an empirical evaluation demonstrating that the new bounds can improve on the work of Masegosa et al. [2020]. Both the parametric form of the Chebyshev-Cantelli inequality and the PAC-Bayes-Bennett inequality may be of independent interest for the study of concentration of measure in other domains.

[A Multi-Implicit Neural Representation for Fonts](#)

- Pradyumna Reddy, Zhifei Zhang, Zhaowen Wang, Matthew Fisher, Hailin Jin, Niloy Mitra
- abstract@[open-review](#): Fonts are ubiquitous across documents and come in a variety of styles. They are either represented in a native vector format or rasterized to produce fixed resolution images. In the first case, the non-standard representation prevents benefiting from latest network architectures for neural representations; while, in the latter case, the rasterized representation, when encoded via networks, results in loss of data fidelity, as font-specific discontinuities like edges and corners are difficult to represent using neural networks. Based on the observation that complex fonts can be represented by a superposition of a set of simpler occupancy functions, we introduce multi-implicits to represent fonts as a permutation-invariant set of learned implicit functions, without losing features (e.g., edges and corners). However, while multi-implicits locally preserve font features, obtaining supervision in the form of ground truth multi-channel signals is a problem in itself. Instead, we propose how to train such a representation with only local supervision, while the proposed neural architecture directly finds globally consistent multi-implicits for font families. We extensively evaluate the proposed representation for various tasks including reconstruction, interpolation, and synthesis to demonstrate clear advantages with existing alternatives. Additionally, the representation naturally enables glyph completion, wherein a single characteristic font is used to synthesize a whole font family in the target style.

[OctField: Hierarchical Implicit Functions for 3D Modeling](#)

- Jia-Heng Tang, Weikai Chen, jie Yang, Bo Wang, Songrun Liu, Bo Yang, Lin Gao
- abstract@[open-review](#): Recent advances in localized implicit functions have enabled neural implicit representation to be scalable to large scenes. However, the regular subdivision of 3D space employed by these approaches fails to take into account the sparsity of the surface occupancy and the varying granularities of geometric details. As a result, its memory footprint grows cubically with the input volume, leading to a prohibitive computational cost even at a moderately dense decomposition. In this work, we present a learnable hierarchical implicit representation for 3D surfaces, coded OctField, that allows high-precision encoding of intricate surfaces with low memory and computational budget. The key to our approach is an adaptive decomposition of 3D scenes that only distributes local implicit functions around the surface of interest. We achieve this goal by introducing a hierarchical octree structure to adaptively subdivide the 3D space according to the surface occupancy and the richness of part geometry. As octree is discrete and non-differentiable, we further propose a novel hierarchical network that models the subdivision of octree cells as a probabilistic process and recursively encodes and decodes both octree structure and surface geometry in a differentiable manner. We demonstrate the value of OctField for a range of shape modeling and reconstruction tasks, showing superiority over alternative approaches.

[The Inductive Bias of Quantum Kernels](#)

- Jonas Kübler, Simon Buchholz, Bernhard Schölkopf
- abstract@[open-review](#): It has been hypothesized that quantum computers may lend themselves well to applications in machine learning. In the present work, we analyze function classes defined via quantum kernels. Quantum computers offer the possibility to efficiently compute inner products of exponentially large density operators that are classically hard to compute. However, having an exponentially large feature space renders the problem of generalization hard. Furthermore, being able to evaluate inner products in high dimensional spaces efficiently by itself does not guarantee a quantum advantage, as already classically tractable kernels can correspond to high- or infinite-dimensional reproducing kernel Hilbert spaces (RKHS). We analyze the spectral properties of quantum kernels and find that we can expect an advantage if their RKHS is low dimensional and contains functions that are hard to compute classically. If the target function is known to lie in this class, this implies a quantum advantage, as the quantum computer can encode this inductive bias, whereas there is no classically efficient way to constrain the function class in the same way. However, we show that finding suitable quantum kernels is not easy because the kernel evaluation might require exponentially many measurements. In conclusion, our message is a somewhat sobering one: we conjecture that quantum machine learning models can offer speed-ups only if we manage to encode knowledge about the problem at hand into quantum circuits, while encoding the same bias into a classical model would be hard. These situations may plausibly occur when learning on data generated by a quantum process, however, they appear to be harder to come by for classical datasets.

[An Exponential Improvement on the Memorization Capacity of Deep Threshold Networks](#)

- Shashank Rajput, Kartik Sreenivasan, Dimitris Papailiopoulos, Amin Karbasi
- abstract@[open-review](#): It is well known that modern deep neural networks are powerful enough to memorize datasets even when the labels have been randomized. Recently, Vershynin(2020) settled a long standing question by Baum(1988), proving that deep threshold networks can memorize n points in d dimensions using $\tilde{O}(e^{1/\delta^2} + \sqrt{n})$ neurons and $\tilde{O}(e^{1/\delta^2}(d + \sqrt{n}) + n)$ weights, where δ is the minimum distance between the points. In this work, we improve the dependence on δ from exponential to almost linear, proving that $\tilde{O}(\frac{1}{\delta} + \sqrt{n})$ neurons and $\tilde{O}(\frac{d}{\delta} + n)$ weights are sufficient. Our construction uses Gaussian random weights only in the first layer, while all the subsequent layers use binary or integer weights. We also prove new lower bounds by connecting memorization in neural networks to the purely geometric problem of separating n points on a sphere using hyperplanes.

[Pretraining Representations for Data-Efficient Reinforcement Learning](#)

- Max Schwarzer, Nitarshan Rajkumar, Michael Noukhovitch, Ankesh Anand, Laurent Charlin, R Devon Hjelm, Philip Bachman, Aaron C. Courville
- abstract@[open-review](#): Data efficiency is a key challenge for deep reinforcement learning. We address this problem by using unlabeled data to pretrain an encoder which is then finetuned on a small amount of task-specific data. To encourage learning representations which capture diverse aspects of the underlying MDP, we employ a combination of latent dynamics modelling and unsupervised goal-conditioned RL. When limited to 100k steps of interaction on Atari games (equivalent to two hours of human experience), our approach significantly surpasses prior work combining offline representation pretraining with task-specific finetuning, and compares favourably with other pretraining methods that require orders of magnitude more data. Our approach shows particular promise when combined with larger models as well as more diverse, task-aligned observational data -- approaching human-level performance and data-efficiency on Atari in our best setting.

[Universal Approximation Using Well-Conditioned Normalizing Flows](#)

- Holden Lee, Chirag Pabbaraju, Anish Prasad Sevukari, Andrej Risteski
- abstract@[open-review](#): Normalizing flows are a widely used class of latent-variable generative models with a tractable likelihood. Affine-coupling models [Dinh et al., 2014, 2016] are a particularly common type of normalizing flows, for which the Jacobian of the latent-to-observable-variable transformation

is triangular, allowing the likelihood to be computed in linear time. Despite the widespread usage of affine couplings, the special structure of the architecture makes understanding their representational power challenging. The question of universal approximation was only recently resolved by three parallel papers [Huang et al., 2020, Zhang et al., 2020, Koehler et al., 2020] – who showed reasonably regular distributions can be approximated arbitrarily well using affine couplings – albeit with networks with a nearly-singular Jacobian. As ill-conditioned Jacobians are an obstacle for likelihood-based training, the fundamental question remains: which distributions can be approximated using well-conditioned affine coupling flows? In this paper, we show that any log-concave distribution can be approximated using well-conditioned affine-coupling flows. In terms of proof techniques, we uncover and leverage deep connections between affine coupling architectures, underdamped Langevin dynamics (a stochastic differential equation often used to sample from Gibbs measures) and Hénon maps (a structured dynamical system that appears in the study of symplectic diffeomorphisms). In terms of informing practice, we approximate a padded version of the input distribution with iid Gaussians – a strategy which Koehler et al. [2020] empirically observed to result in better-conditioned flows, but had hitherto no theoretical grounding. Our proof can thus be seen as providing theoretical evidence for the benefits of Gaussian padding when training normalizing flows.

[On the Validity of Modeling SGD with Stochastic Differential Equations \(SDEs\)](#)

- Zhiyuan Li, Sadhika Malladi, Sanjeev Arora
- abstract@[open-review](#): It is generally recognized that finite learning rate (LR), in contrast to infinitesimal LR, is important for good generalization in real-life deep nets. Most attempted explanations propose approximating finite-LR SGD with Itô Stochastic Differential Equations (SDEs), but formal justification for this approximation (e.g., Li et al., 2019) only applies to SGD with tiny LR. Experimental verification of the approximation appears computationally infeasible. The current paper clarifies the picture with the following contributions: (a) An efficient simulation algorithm SVAG that provably converges to the conventionally used Itô SDE approximation. (b) A theoretically motivated testable necessary condition for the SDE approximation and its most famous implication, the linear scaling rule (Goyal et al., 2017), to hold.(c) Experiments using this simulation to demonstrate that the previously proposed SDE approximation can meaningfully capture the training and generalization properties of common deep nets.

[Proportional Participatory Budgeting with Additive Utilities](#)

- Grzegorz Pierczyński, Piotr Skowron, Dominik Peters
- abstract@[open-review](#): We study voting rules for participatory budgeting, where a group of voters collectively decides which projects should be funded using a common budget. We allow the projects to have arbitrary costs, and the voters to have arbitrary additive valuations over the projects. We formulate two axioms that guarantee proportional representation to groups of voters with common interests. To the best of our knowledge, all known rules for participatory budgeting do not satisfy either of the two axioms; in addition we show that the most prominent proportional rule for committee elections, Proportional Approval Voting, cannot be adapted to arbitrary costs nor to additive valuations so that it would satisfy our axioms of proportionality. We construct a simple and attractive voting rule that satisfies one of our axioms (for arbitrary costs and arbitrary additive valuations), and that can be evaluated in polynomial time. We prove that our other stronger axiom is also satisfiable, though by a computationally more expensive and less natural voting rule.

[Disentangling the Roles of Curation, Data-Augmentation and the Prior in the Cold Posterior Effect](#)

- Lorenzo Noci, Kevin Roth, Gregor Bachmann, Sebastian Nowozin, Thomas Hofmann
- abstract@[open-review](#): The “cold posterior effect” (CPE) in Bayesian deep learning describes the disturbing observation that the predictive performance of Bayesian neural networks can be significantly improved if the Bayes posterior is artificially sharpened using a temperature parameter $T < 1$. The CPE is problematic in theory and practice and since the effect was identified many researchers have proposed hypotheses to explain the phenomenon. However, despite this intensive research effort the effect remains poorly understood. In this work we provide novel and nuanced evidence relevant to existing explanations for the cold posterior effect, disentangling three hypotheses: 1. The dataset curation hypothesis of Aitchison (2020): we show empirically that the CPE does not arise in a real curated data set but can be produced in a controlled experiment with varying curation strength. 2. The data augmentation hypothesis of Izmailov et al. (2021) and Fortuin et al. (2021): we show empirically that data augmentation is sufficient but not necessary for the CPE to be present. 3. The bad prior hypothesis of Wenzel et al. (2020): we use a simple experiment evaluating the relative importance of the prior and the likelihood, strongly linking the CPE to the prior. Our results demonstrate how the CPE can arise in isolation from synthetic curation, data augmentation, and bad priors. Cold posteriors observed “in the wild” are therefore unlikely to arise from a single simple cause; as a result, we do not expect a simple “fix” for cold posteriors.

[Sanity Checks for Lottery Tickets: Does Your Winning Ticket Really Win the Jackpot?](#)

- Xiaolong Ma, Geng Yuan, Xuan Shen, Tianlong Chen, Xuxi Chen, Xiaohan Chen, Ning Liu, Minghai Qin, Sijia Liu, Zhangyang Wang, Yanzhi Wang
- abstract@[open-review](#): There have been long-standing controversies and inconsistencies over the experiment setup and criteria for identifying the “winning ticket” in literature. To reconcile such, we revisit the definition of lottery ticket hypothesis, with comprehensive and more rigorous conditions. Under our new definition, we show concrete evidence to clarify whether the winning ticket exists across the major DNN architectures and/or applications. Through extensive experiments, we perform quantitative analysis on the correlations between winning tickets and various experimental factors, and empirically study the patterns of our observations. We find that the key training hyperparameters, such as learning rate and training epochs, as well as the architecture characteristics such as capacities and residual connections, are all highly correlated with whether and when the winning tickets can be identified. Based on our analysis, we summarize a guideline for parameter settings in regards of specific architecture characteristics, which we hope to catalyze the research progress on the topic of lottery ticket hypothesis. Our codes are publicly available at: <https://github.com/boone891214/sanity-check-LTH>.

[Collaborative Causal Discovery with Atomic Interventions](#)

- Raghavendra Addanki, Shiva Kasiviswanathan
- abstract@[open-review](#): We introduce a new Collaborative Causal Discovery problem, through which we model a common scenario in which we have multiple independent entities each with their own causal graph, and the goal is to simultaneously learn all these causal graphs. We study this problem without the causal sufficiency assumption, using Maximal Ancestral Graphs (MAG) to model the causal graphs, and assuming that we have the ability to actively perform independent single vertex (or atomic) interventions on the entities. If the $\$M\$$ underlying (unknown) causal graphs of the entities satisfy a natural notion of clustering, we give algorithms that leverage this property and recovers all the causal graphs using roughly logarithmic in $\$M\$$ number of atomic interventions per entity. These are significantly fewer than $\$n\$$ atomic interventions per entity required to learn each causal graph separately, where $\$n\$$ is the number of observable nodes in the causal graph. We complement our results with a lower bound and discuss various extensions of our collaborative setting.

[Towards optimally abstaining from prediction with OOD test examples](#)

- Adam Kalai, Varun Kanade
- abstract@[open-review](#): A common challenge across all areas of machine learning is that training data is not distributed like test data, due to natural shifts or adversarial examples; such examples are referred to as out-of-distribution (OOD) test examples. We consider a model where one may abstain from predicting, at a fixed cost. In particular, our transductive abstention algorithm takes labeled training examples and unlabeled test examples as input, and provides predictions with optimal prediction loss guarantees. The loss bounds match standard generalization bounds when test examples are i.i.d. from the

training distribution, but add an additional term that is the cost of abstaining times the statistical distance between the train and test distribution (or the fraction of adversarial examples). For linear regression, we give a polynomial-time algorithm based on Celis-Dennis-Tapia optimization algorithms. For binary classification, we show how to efficiently implement it using a proper agnostic learner (i.e., an Empirical Risk Minimizer) for the class of interest. Our work builds on recent work of Goldwasser, Kalai, and Montasser (2020) who gave error and abstention guarantees for transductive binary classification.

[TokenLearner: Adaptive Space-Time Tokenization for Videos](#)

- Michael Ryoo, AJ Piergiovanni, Anurag Arnab, Mostafa Dehghani, Anelia Angelova
- abstract@[open-review](#): In this paper, we introduce a novel visual representation learning which relies on a handful of adaptively learned tokens, and which is applicable to both image and video understanding tasks. Instead of relying on hand-designed splitting strategies to obtain visual tokens and processing a large number of densely sampled patches for attention, our approach learns to mine important tokens in visual data. This results in efficiently and effectively finding a few important visual tokens and enables modeling of pairwise attention between such tokens, over a longer temporal horizon for videos, or the spatial content in image frames. Our experiments demonstrate strong performance on several challenging benchmarks for video recognition tasks. Importantly, due to our tokens being adaptive, we accomplish competitive results at significantly reduced computational cost. We establish new state-of-the-arts on multiple video datasets, including Kinetics-400, Kinetics-600, Charades, and AViD.

[Learning in Multi-Stage Decentralized Matching Markets](#)

- Xiaowu Dai, Michael Jordan
- abstract@[open-review](#): Matching markets are often organized in a multi-stage and decentralized manner. Moreover, participants in real-world matching markets often have uncertain preferences. This article develops a framework for learning optimal strategies in such settings, based on a nonparametric statistical approach and variational analysis. We propose an efficient algorithm, built upon concepts of "lower uncertainty bound" and "calibrated decentralized matching," for maximizing the participants' expected payoff. We show that there exists a welfare-versus-fairness trade-off that is characterized by the uncertainty level of acceptance. Participants will strategically act in favor of a low uncertainty level to reduce competition and increase expected payoff. We prove that participants can be better off with multi-stage matching compared to single-stage matching. We demonstrate aspects of the theoretical predictions through simulations and an experiment using real data from college admissions.

[Non-asymptotic convergence bounds for Wasserstein approximation using point clouds](#)

- Quentin Mérigot, Filippo Santambrogio, Clément SARAZIN
- abstract@[open-review](#): Several issues in machine learning and inverse problems require to generate discrete data, as if sampled from a model probability distribution. A common way to do so relies on the construction of a uniform probability distribution over a set of N points which minimizes the Wasserstein distance to the model distribution. This minimization problem, where the unknowns are the positions of the atoms, is non-convex. Yet, in most cases, a suitably adjusted version of Lloyd's algorithm in which Voronoi cells are replaced by Power cells, leads to configurations with small Wasserstein error. This is surprising because, again, of the non-convex nature of the problem, which moreover admits spurious critical points. We provide explicit upper bounds for the convergence speed of this Lloyd-type algorithm, starting from a cloud of points sufficiently far from each other. This already works after one step of the iteration procedure, and similar bounds can be deduced, for the corresponding gradient descent. These bounds naturally lead to a sort of Poliak-Łojasiewicz inequality for the Wasserstein distance cost, with an error term depending on the distances between Dirac masses in the discrete distribution.

[Understanding Interlocking Dynamics of Cooperative Rationalization](#)

- Mo Yu, Yang Zhang, Shiyu Chang, Tommi Jaakkola
- abstract@[open-review](#): Selective rationalization explains the prediction of complex neural networks by finding a small subset of the input that is sufficient to predict the neural model output. The selection mechanism is commonly integrated into the model itself by specifying a two-component cascaded system consisting of a rationale generator, which makes a binary selection of the input features (which is the rationale), and a predictor, which predicts the output based only on the selected features. The components are trained jointly to optimize prediction performance. In this paper, we reveal a major problem with such cooperative rationalization paradigm --- model interlocking. Inter-locking arises when the predictor overfits to the features selected by the generator thus reinforcing the generator's selection even if the selected rationales are sub-optimal. The fundamental cause of the interlocking problem is that the rationalization objective to be minimized is concave with respect to the generator's selection policy. We propose a new rationalization framework, called A2R, which introduces a third component into the architecture, a predictor driven by soft attention as opposed to selection. The generator now realizes both soft and hard attention over the features and these are fed into the two different predictors. While the generator still seeks to support the original predictor performance, it also minimizes a gap between the two predictors. As we will show theoretically, since the attention-based predictor exhibits a better convexity property, A2R can overcome the concavity barrier. Our experiments on two synthetic benchmarks and two real datasets demonstrate that A2R can significantly alleviate the interlock problem and find explanations that better align with human judgments.

[Adversarial Robustness without Adversarial Training: A Teacher-Guided Curriculum Learning Approach](#)

- Anindya Sarkar, Anirban Sarkar, Sowrya Gali, Vineeth N Balasubramanian
- abstract@[open-review](#): Current SOTA adversarially robust models are mostly based on adversarial training (AT) and differ only by some regularizers either at inner maximization or outer minimization steps. Being repetitive in nature during the inner maximization step, they take a huge time to train. We propose a non-iterative method that enforces the following ideas during training. Attribution maps are more aligned to the actual object in the image for adversarially robust models compared to naturally trained models. Also, the allowed set of pixels to perturb an image (that changes model decision) should be restricted to the object pixels only, which reduces the attack strength by limiting the attack space. Our method achieves significant performance gains with a little extra effort (10-20%) over existing AT models and outperforms all other methods in terms of adversarial as well as natural accuracy. We have performed extensive experimentation with CIFAR-10, CIFAR-100, and TinyImageNet datasets and reported results against many popular strong adversarial attacks to prove the effectiveness of our method.

[Tactical Optimism and Pessimism for Deep Reinforcement Learning](#)

- Ted Moskovitz, Jack Parker-Holder, Aldo Pacchiano, Michael Arbel, Michael Jordan
- abstract@[open-review](#): In recent years, deep off-policy actor-critic algorithms have become a dominant approach to reinforcement learning for continuous control. One of the primary drivers of this improved performance is the use of pessimistic value updates to address function approximation errors, which previously led to disappointing performance. However, a direct consequence of pessimism is reduced exploration, running counter to theoretical support for the efficacy of optimism in the face of uncertainty. So which approach is best? In this work, we show that the most effective degree of optimism can vary both across tasks and over the course of learning. Inspired by this insight, we introduce a novel deep actor-critic framework, Tactical Optimistic and Pessimistic (TOP) estimation, which switches between optimistic and pessimistic value learning online. This is achieved by formulating the selection as a multi-arm bandit problem. We show in a series of continuous control tasks that TOP outperforms existing methods which rely on a fixed degree of optimism, setting a new state of the art in challenging pixel-based environments. Since our changes are simple to implement, we believe these insights can easily be incorporated into a multitude of off-policy algorithms.

[Towards Hyperparameter-free Policy Selection for Offline Reinforcement Learning](#)

- Siyuan Zhang, Nan Jiang
- abstract@[open-review](#): How to select between policies and value functions produced by different training algorithms in offline reinforcement learning (RL)---which is crucial for hyperparameter tuning---is an important open question. Existing approaches based on off-policy evaluation (OPE) often require additional function approximation and hence hyperparameters, creating a chicken-and-egg situation. In this paper, we design hyperparameter-free algorithms for policy selection based on BVFT [XJ21], a recent theoretical advance in value-function selection, and demonstrate their effectiveness in discrete-action benchmarks such as Atari. To address performance degradation due to poor critics in continuous-action domains, we further combine BVFT with OPE to get the best of both worlds, and obtain a hyperparameter-tuning method for Q\$-function based OPE with theoretical guarantees as a side product.

[FjORD: Fair and Accurate Federated Learning under heterogeneous targets with Ordered Dropout](#)

- Samuel Horváth, Stefanos Laskaridis, Mario Almeida, Ilias Leontiadis, Stylianos Venieris, Nicholas Lane
- abstract@[open-review](#): Federated Learning (FL) has been gaining significant traction across different ML tasks, ranging from vision to keyboard predictions. In large-scale deployments, client heterogeneity is a fact and constitutes a primary problem for fairness, training performance and accuracy. Although significant efforts have been made into tackling statistical data heterogeneity, the diversity in the processing capabilities and network bandwidth of clients, termed system heterogeneity, has remained largely unexplored. Current solutions either disregard a large portion of available devices or set a uniform limit on the model's capacity, restricted by the least capable participants. In this work, we introduce Ordered Dropout, a mechanism that achieves an ordered, nested representation of knowledge in Neural Networks and enables the extraction of lower footprint submodels without the need for retraining. We further show that for linear maps our Ordered Dropout is equivalent to SVD. We employ this technique, along with a self-distillation methodology, in the realm of FL in a framework called FjORD. FjORD alleviates the problem of client system heterogeneity by tailoring the model width to the client's capabilities. Extensive evaluation on both CNNs and RNNs across diverse modalities shows that FjORD consistently leads to significant performance gains over state-of-the-art baselines while maintaining its nested structure.

[Optimal Uniform OPE and Model-based Offline Reinforcement Learning in Time-Homogeneous, Reward-Free and Task-Agnostic Settings](#)

- Ming Yin, Yu-Xiang Wang
- abstract@[open-review](#): This work studies the statistical limits of uniform convergence for offline policy evaluation (OPE) problems with model-based methods (for episodic MDP) and provides a unified framework towards optimal learning for several well-motivated offline tasks. Uniform OPE $\sup_{\Pi} \mathbb{P}[\hat{Q}^{\pi} - Q^{\pi}] \leq \epsilon$ is a stronger measure than the point-wise OPE and ensures offline learning when Π contains all policies (the global class). In this paper, we establish an $\Omega(\mathcal{H}^2 S/d_m \epsilon^2)$ lower bound (over model-based family) for the global uniform OPE and our main result establishes an upper bound of $\tilde{O}(\mathcal{H}^2/d_m \epsilon^2)$ for the local uniform convergence that applies to all near-empirically optimal policies for the MDPs with stationary transition. Here d_m is the minimal marginal state-action probability. Critically, the highlight in achieving the optimal rate $\tilde{O}(\mathcal{H}^2/d_m \epsilon^2)$ is our design of singleton absorbing MDP, which is a new sharp analysis tool that works with the model-based approach. We generalize such a model-based framework to the new settings: offline task-agnostic and the offline reward-free with optimal complexity $\tilde{O}(\mathcal{H}^2 \log(K)/d_m \epsilon^2)$ (K is the number of tasks) and $\tilde{O}(\mathcal{H}^2 S/d_m \epsilon^2)$ respectively. These results provide a unified solution for simultaneously solving different offline RL problems.

[MixSeq: Connecting Macroscopic Time Series Forecasting with Microscopic Time Series Data](#)

- Zhibo Zhu, Ziqi Liu, Ge Jin, Zhiqiang Zhang, Lei Chen, Jun Zhou, Jianyong Zhou
- abstract@[open-review](#): Time series forecasting is widely used in business intelligence, e.g., forecast stock market price, sales, and help the analysis of data trend. Most time series of interest are macroscopic time series that are aggregated from microscopic data. However, instead of directly modeling the macroscopic time series, rare literature studied the forecasting of macroscopic time series by leveraging data on the microscopic level. In this paper, we assume that the microscopic time series follow some unknown mixture probabilistic distributions. We theoretically show that as we identify the ground truth latent mixture components, the estimation of time series from each component could be improved because of lower variance, thus benefitting the estimation of macroscopic time series as well. Inspired by the power of Seq2seq and its variants on the modeling of time series data, we propose Mixture of Seq2seq (MixSeq), an end2end mixture model to cluster microscopic time series, where all the components come from a family of Seq2seq models parameterized by different parameters. Extensive experiments on both synthetic and real-world data show the superiority of our approach.

[Pareto Domain Adaptation](#)

- fangrui lv, Jian Liang, Kaixiong Gong, Shuang Li, Chi Harold Liu, Han Li, Di Liu, Guoren Wang
- abstract@[open-review](#): Domain adaptation (DA) attempts to transfer the knowledge from a labeled source domain to an unlabeled target domain that follows different distribution from the source. To achieve this, DA methods include a source classification objective to extract the source knowledge and a domain alignment objective to diminish the domain shift, ensuring knowledge transfer. Typically, former DA methods adopt some weight hyper-parameters to linearly combine the training objectives to form an overall objective. However, the gradient directions of these objectives may conflict with each other due to domain shift. Under such circumstances, the linear optimization scheme might decrease the overall objective value at the expense of damaging one of the training objectives, leading to restricted solutions. In this paper, we rethink the optimization scheme for DA from a gradient-based perspective. We propose a Pareto Domain Adaptation (ParetoDA) approach to control the overall optimization direction, aiming to cooperatively optimize all training objectives. Specifically, to reach a desirable solution on the target domain, we design a surrogate loss mimicking target classification. To improve target-prediction accuracy to support the mimicking, we propose a target-prediction refining mechanism which exploits domain labels via Bayes' theorem. On the other hand, since prior knowledge of weighting schemes for objectives is often unavailable to guide optimization to approach the optimal solution on the target domain, we propose a dynamic preference mechanism to dynamically guide our cooperative optimization by the gradient of the surrogate loss on a held-out unlabeled target dataset. Our theoretical analyses show that the held-out data can guide but will not be over-fitted by the optimization. Extensive experiments on image classification and semantic segmentation benchmarks demonstrate the effectiveness of ParetoDA

[Divergence Frontiers for Generative Models: Sample Complexity, Quantization Effects, and Frontier Integrals](#)

- Lang Liu, Krishna Pillutla, Sean Welleck, Sewoong Oh, Yejin Choi, Zaid Harchaoui
- abstract@[open-review](#): The spectacular success of deep generative models calls for quantitative tools to measure their statistical performance. Divergence frontiers have recently been proposed as an evaluation framework for generative models, due to their ability to measure the quality-diversity trade-off inherent to deep generative modeling. We establish non-asymptotic bounds on the sample complexity of divergence frontiers. We also introduce frontier integrals which provide summary statistics of divergence frontiers. We show how smoothed estimators such as Good-Turing or Krichevsky-Trofimov can overcome the missing mass problem and lead to faster rates of convergence. We illustrate the theoretical results with numerical examples from natural language processing and computer vision.

[Consistency Regularization for Variational Auto-Encoders](#)

- Samarth Sinha, Adji Bouso Dieng
- abstract@[open-review](#): Variational Auto-Encoders (VAEs) are a powerful approach to unsupervised learning. They enable scalable approximate posterior inference in latent-variable models using variational inference. A VAE posits a variational family parameterized by a deep neural network---called an encoder---that takes data as input. This encoder is shared across all the observations, which amortizes the cost of inference. However the encoder of a VAE has the undesirable property that it maps a given observation and a semantics-preserving transformation of it to different latent representations. This "inconsistency" of the encoder lowers the quality of the learned representations, especially for downstream tasks, and also negatively affects generalization. In this paper, we propose a regularization method to enforce consistency in VAEs. The idea is to minimize the Kullback-Leibler (KL) divergence between the variational distribution when conditioning on the observation and the variational distribution when conditioning on a random semantics-preserving transformation of this observation. This regularization is applicable to any VAE. In our experiments we apply it to four different VAE variants on several benchmark datasets and found it always improves the quality of the learned representations but also leads to better generalization. In particular, when applied to the Nouveau VAE (NVAE), our regularization method yields state-of-the-art performance on MNIST, CIFAR-10, and CELEBA. We also applied our method to 3D data and found it learns representations of superior quality as measured by accuracy on a downstream classification task. Finally, we show our method can even outperform the triplet loss, an advanced and popular contrastive learning-based method for representation learning.

[Score-based Generative Neural Networks for Large-Scale Optimal Transport](#)

- Max Daniels, Tyler Maunu, Paul Hand
- abstract@[open-review](#): We consider the fundamental problem of sampling the optimal transport coupling between given source and target distributions. In certain cases, the optimal transport plan takes the form of a one-to-one mapping from the source support to the target support, but learning or even approximating such a map is computationally challenging for large and high-dimensional datasets due to the high cost of linear programming routines and an intrinsic curse of dimensionality. We study instead the Sinkhorn problem, a regularized form of optimal transport whose solutions are couplings between the source and the target distribution. We introduce a novel framework for learning the Sinkhorn coupling between two distributions in the form of a score-based generative model. Conditioned on source data, our procedure iterates Langevin Dynamics to sample target data according to the regularized optimal coupling. Key to this approach is a neural network parametrization of the Sinkhorn problem, and we prove convergence of gradient descent with respect to network parameters in this formulation. We demonstrate its empirical success on a variety of large scale optimal transport tasks.

[Interactive Label Cleaning with Example-based Explanations](#)

- Stefano Teso, Andrea Bontempelli, Fausto Giunchiglia, Andrea Passerini
- abstract@[open-review](#): We tackle sequential learning under label noise in applications where a human supervisor can be queried to relabel suspicious examples. Existing approaches are flawed, in that they only relabel incoming examples that look "suspicious" to the model. As a consequence, those mislabeled examples that elude (or don't undergo) this cleaning step end up tainting the training data and the model with no further chance of being cleaned. We propose CINCER, a novel approach that cleans both new and past data by identifying \emph{pairs of mutually incompatible examples}. Whenever it detects a suspicious example, CINCER identifies a counter-example in the training set that - according to the model - is maximally incompatible with the suspicious example, and asks the annotator to relabel either or both examples, resolving this possible inconsistency. The counter-examples are chosen to be maximally incompatible, so to serve as \emph{explanations} of the model's suspicion, and highly influential, so to convey as much information as possible if relabeled. CINCER achieves this by leveraging an efficient and robust approximation of influence functions based on the Fisher information matrix (FIM). Our extensive empirical evaluation shows that clarifying the reasons behind the model's suspicions by cleaning the counter-examples helps in acquiring substantially better data and models, especially when paired with our FIM approximation.

[Gradient Descent on Two-layer Nets: Margin Maximization and Simplicity Bias](#)

- Kaifeng Lyu, Zhiyuan Li, Runzhe Wang, Sanjeev Arora
- abstract@[open-review](#): The generalization mystery of overparametrized deep nets has motivated efforts to understand how gradient descent (GD) converges to low-loss solutions that generalize well. Real-life neural networks are initialized from small random values and trained with cross-entropy loss for classification (unlike the "lazy" or "NTK" regime of training where analysis was more successful), and a recent sequence of results (Lyu and Li, 2020; Chizat and Bach, 2020; Ji and Telgarsky, 2020) provide theoretical evidence that GD may converge to the "max-margin" solution with zero loss, which presumably generalizes well. However, the global optimality of margin is proved only in some settings where neural nets are infinitely or exponentially wide. The current paper is able to establish this global optimality for two-layer Leaky ReLU nets trained with gradient flow on linearly separable and symmetric data, regardless of the width. The analysis also gives some theoretical justification for recent empirical findings (Kalimeris et al., 2019) on the so-called simplicity bias of GD towards linear or other "simple" classes of solutions, especially early in training. On the pessimistic side, the paper suggests that such results are fragile. A simple data manipulation can make gradient flow converge to a linear classifier with suboptimal margin.

[Glance-and-Gaze Vision Transformer](#)

- Qihang Yu, Yingda Xia, Yutong Bai, Yongyi Lu, Alan L. Yuille, Wei Shen
- abstract@[open-review](#): Recently, there emerges a series of vision Transformers, which show superior performance with a more compact model size than conventional convolutional neural networks, thanks to the strong ability of Transformers to model long-range dependencies. However, the advantages of vision Transformers also come with a price: Self-attention, the core part of Transformer, has a quadratic complexity to the input sequence length. This leads to a dramatic increase of computation and memory cost with the increase of sequence length, thus introducing difficulties when applying Transformers to the vision tasks that require dense predictions based on high-resolution feature maps. In this paper, we propose a new vision Transformer, named Glance-and-Gaze Transformer (GG-Transformer), to address the aforementioned issues. It is motivated by the Glance and Gaze behavior of human beings when recognizing objects in natural scenes, with the ability to efficiently model both long-range dependencies and local context. In GG-Transformer, the Glance and Gaze behavior is realized by two parallel branches: The Glance branch is achieved by performing self-attention on the adaptively-dilated partitions of the input, which leads to a linear complexity while still enjoying a global receptive field; The Gaze branch is implemented by a simple depth-wise convolutional layer, which compensates local image context to the features obtained by the Glance mechanism. We empirically demonstrate our method achieves consistently superior performance over previous state-of-the-art Transformers on various vision tasks and benchmarks.

[Stochastic \$L^{\text{natural}}\$ -convex Function Minimization](#)

- Haixiang Zhang, Zeyu Zheng, Javad Lavaei
- abstract@[open-review](#): We study an extension of the stochastic submodular minimization problem, namely, the stochastic L^{natural} -convex minimization problem. We develop the first polynomial-time algorithms that return a near-optimal solution with high probability. We design a novel truncation operation to further reduce the computational complexity of the proposed algorithms. When applied to a stochastic submodular function, the computational complexity of the proposed algorithms is lower than that of the existing stochastic submodular minimization algorithms. In addition, we provide a strongly polynomial approximate algorithm. The algorithm execution also does not require any prior knowledge about the objective function except the L^{natural} -convexity. A lower bound on the computational complexity that is required to achieve a high probability error bound is also derived. Numerical experiments are implemented to demonstrate the efficiency of our theoretical findings.

[Self-Supervised GANs with Label Augmentation](#)

- Liang Hou, Huawei Shen, Qi Cao, Xueqi Cheng
- abstract@[open-review](#): Recently, transformation-based self-supervised learning has been applied to generative adversarial networks (GANs) to mitigate catastrophic forgetting in the discriminator by introducing a stationary learning environment. However, the separate self-supervised tasks in existing self-supervised GANs cause a goal inconsistent with generative modeling due to the fact that their self-supervised classifiers are agnostic to the generator distribution. To address this problem, we propose a novel self-supervised GAN that unifies the GAN task with the self-supervised task by augmenting the GAN labels (real or fake) via self-supervision of data transformation. Specifically, the original discriminator and self-supervised classifier are unified into a label-augmented discriminator that predicts the augmented labels to be aware of both the generator distribution and the data distribution under every transformation, and then provide the discrepancy between them to optimize the generator. Theoretically, we prove that the optimal generator could converge to replicate the real data distribution. Empirically, we show that the proposed method significantly outperforms previous self-supervised and data augmentation GANs on both generative modeling and representation learning across benchmark datasets.

[Shape As Points: A Differentiable Poisson Solver](#)

- Songyou Peng, Chiyu Jiang, Yiyi Liao, Michael Niemeyer, Marc Pollefeys, Andreas Geiger
- abstract@[open-review](#): In recent years, neural implicit representations gained popularity in 3D reconstruction due to their expressiveness and flexibility. However, the implicit nature of neural implicit representations results in slow inference times and requires careful initialization. In this paper, we revisit the classic yet ubiquitous point cloud representation and introduce a differentiable point-to-mesh layer using a differentiable formulation of Poisson Surface Reconstruction (PSR) which allows for a GPU-accelerated fast solution of the indicator function given an oriented point cloud. The differentiable PSR layer allows us to efficiently and differentiably bridge the explicit 3D point representation with the 3D mesh via the implicit indicator field, enabling end-to-end optimization of surface reconstruction metrics such as Chamfer distance. This duality between points and meshes hence allows us to represent shapes as oriented point clouds, which are explicit, lightweight and expressive. Compared to neural implicit representations, our Shape-As-Points (SAP) model is more interpretable, lightweight, and accelerates inference time by one order of magnitude. Compared to other explicit representations such as points, patches, and meshes, SAP produces topology-agnostic, watertight manifold surfaces. We demonstrate the effectiveness of SAP on the task of surface reconstruction from unoriented point clouds and learning-based reconstruction.

[Outcome-Driven Reinforcement Learning via Variational Inference](#)

- Tim G. J. Rudner, Vitchyr Pong, Rowan McAllister, Yarin Gal, Sergey Levine
- abstract@[open-review](#): While reinforcement learning algorithms provide automated acquisition of optimal policies, practical application of such methods requires a number of design decisions, such as manually designing reward functions that not only define the task, but also provide sufficient shaping to accomplish it. In this paper, we view reinforcement learning as inferring policies that achieve desired outcomes, rather than as a problem of maximizing rewards. To solve this inference problem, we establish a novel variational inference formulation that allows us to derive a well-shaped reward function which can be learned directly from environment interactions. From the corresponding variational objective, we also derive a new probabilistic Bellman backup operator and use it to develop an off-policy algorithm to solve goal-directed tasks. We empirically demonstrate that this method eliminates the need to hand-craft reward functions for a suite of diverse manipulation and locomotion tasks and leads to effective goal-directed behaviors.

[Drawing Robust Scratch Tickets: Subnetworks with Inborn Robustness Are Found within Randomly Initialized Networks](#)

- Yonggan Fu, Qixuan Yu, Yang Zhang, Shang Wu, Xu Ouyang, David Cox, Yingyan Lin
- abstract@[open-review](#): Deep Neural Networks (DNNs) are known to be vulnerable to adversarial attacks, i.e., an imperceptible perturbation to the input can mislead DNNs trained on clean images into making erroneous predictions. To tackle this, adversarial training is currently the most effective defense method, by augmenting the training set with adversarial samples generated on the fly. Interestingly, we discover for the first time that there exist subnetworks with inborn robustness, matching or surpassing the robust accuracy of the adversarially trained networks with comparable model sizes, within randomly initialized networks without any model training}, indicating that adversarial training on model weights is not indispensable towards adversarial robustness. We name such subnetworks Robust Scratch Tickets (RSTs), which are also by nature efficient. Distinct from the popular lottery ticket hypothesis, neither the original dense networks nor the identified RSTs need to be trained. To validate and understand this fascinating finding, we further conduct extensive experiments to study the existence and properties of RSTs under different models, datasets, sparsity patterns, and attacks, drawing insights regarding the relationship between DNNs' robustness and their initialization/overparameterization. Furthermore, we identify the poor adversarial transferability between RSTs of different sparsity ratios drawn from the same randomly initialized dense network, and propose a Random RST Switch (R2S) technique, which randomly switches between different RSTs, as a novel defense method built on top of RSTs. We believe our findings about RSTs have opened up a new perspective to study model robustness and extend the lottery ticket hypothesis.

[Rectifying the Shortcut Learning of Background for Few-Shot Learning](#)

- Xu Luo, Longhui Wei, Liangjian Wen, Jinrong Yang, Lingxi Xie, Zenglin Xu, Qi Tian
- abstract@[open-review](#): The category gap between training and evaluation has been characterised as one of the main obstacles to the success of Few-Shot Learning (FSL). In this paper, we for the first time empirically identify image background, common in realistic images, as a shortcut knowledge helpful for in-class classification but ungeneralizable beyond training categories in FSL. A novel framework, COSOC, is designed to tackle this problem by extracting foreground objects in images at both training and evaluation without any extra supervision. Extensive experiments carried on inductive FSL tasks demonstrate the effectiveness of our approaches.

[SEAL: Self-supervised Embodied Active Learning using Exploration and 3D Consistency](#)

- Devendra Singh Chaplot, Murtaza Dalal, Saurabh Gupta, Jitendra Malik, Russ R. Salakhutdinov
- abstract@[open-review](#): In this paper, we explore how we can build upon the data and models of Internet images and use them to adapt to robot vision without requiring any extra labels. We present a framework called Self-supervised Embodied Active Learning (SEAL). It utilizes perception models trained on internet images to learn an active exploration policy. The observations gathered by this exploration policy are labelled using 3D consistency and used to improve the perception model. We build and utilize 3D semantic maps to learn both action and perception in a completely self-supervised manner. The semantic map is used to compute an intrinsic motivation reward for training the exploration policy and for labelling the agent observations using spatio-temporal 3D consistency and label propagation. We demonstrate that the SEAL framework can be used to close the action-perception loop: it improves object detection and instance segmentation performance of a pretrained perception model by just moving around in training environments and the improved perception model can be used to improve Object Goal Navigation.

[Sifting through the noise: Universal first-order methods for stochastic variational inequalities](#)

- Kimon Antonakopoulos, Thomas Pethick, Ali Kavis, Panayotis Mertikopoulos, Volkan Cevher
- abstract@[open-review](#): We examine a flexible algorithmic framework for solving monotone variational inequalities in the presence of randomness and uncertainty. The proposed template encompasses a wide range of popular first-order methods, including dual averaging, dual extrapolation and optimistic gradient algorithms – both adaptive and non-adaptive. Our first result is that the algorithm achieves the optimal rates of convergence for cocoercive problems when the profile of the randomness is known to the optimizer: $\mathcal{O}(1/\sqrt{T})$ for absolute noise profiles, and $\mathcal{O}(1/T)$ for relative ones. Subsequently, we drop all prior knowledge requirements (the absolute/relative variance of the randomness affecting the problem, the operator's cocoercivity constant, etc.), and we analyze an adaptive instance of the method that gracefully interpolates between the above rates – i.e. it

achieves $\mathcal{O}(1/\sqrt{T})$ and $\mathcal{O}(1/T)$ in the absolute and relative cases, respectively. To our knowledge, this is the first universality result of its kind in the literature and, somewhat surprisingly, it shows that an extra-gradient proxy step is not required to achieve optimal rates.

[Accommodating Picky Customers: Regret Bound and Exploration Complexity for Multi-Objective Reinforcement Learning](#)

- Jingfeng Wu, Vladimir Braverman, Lin Yang
- abstract@[open-review](#): In this paper we consider multi-objective reinforcement learning where the objectives are balanced using preferences. In practice, the preferences are often given in an adversarial manner, e.g., customers can be picky in many applications. We formalize this problem as an episodic learning problem on a Markov decision process, where transitions are unknown and a reward function is the inner product of a preference vector with pre-specified multi-objective reward functions. We consider two settings. In the online setting, the agent receives a (adversarial) preference every episode and proposes policies to interact with the environment. We provide a model-based algorithm that achieves a nearly minimax optimal regret bound $\widetilde{\mathcal{O}}(\sqrt{\min\{d,S\} \cdot H^2 K})$, where d is the number of objectives, S is the number of states, A is the number of actions, H is the length of the horizon, and K is the number of episodes. Furthermore, we consider preference-free exploration, i.e., the agent first interacts with the environment without specifying any preference and then is able to accommodate arbitrary preference vector up to ϵ error. Our proposed algorithm is provably efficient with a nearly optimal trajectory complexity $\widetilde{\mathcal{O}}(\min\{d,S\} \cdot H^3 / \epsilon^2)$. This result partly resolves an open problem raised by [jin2020reward](#).

[Exact Privacy Guarantees for Markov Chain Implementations of the Exponential Mechanism with Artificial Atoms](#)

- Jeremy Seeman, Matthew Reimherr, Aleksandra Slavković
- abstract@[open-review](#): Implementations of the exponential mechanism in differential privacy often require sampling from intractable distributions. When approximate procedures like Markov chain Monte Carlo (MCMC) are used, the end result incurs costs to both privacy and accuracy. Existing work has examined these effects asymptotically, but implementable finite sample results are needed in practice so that users can specify privacy budgets in advance and implement samplers with exact privacy guarantees. In this paper, we use tools from ergodic theory and perfect simulation to design exact finite runtime sampling algorithms for the exponential mechanism by introducing an intermediate modified target distribution using artificial atoms. We propose an additional modification of this sampling algorithm that maintains its ϵ -DP guarantee and has improved runtime at the cost of some utility. We then compare these methods in scenarios where we can explicitly calculate a δ cost (as in (ϵ, δ) -DP) incurred when using standard MCMC techniques. Much as there is a well known trade-off between privacy and utility, we demonstrate that there is also a trade-off between privacy guarantees and runtime.

[The Emergence of Objectness: Learning Zero-shot Segmentation from Videos](#)

- Runtao Liu, Zhirong Wu, Stella Yu, Stephen Lin
- abstract@[open-review](#): Humans can easily detect and segment moving objects simply by observing how they move, even without knowledge of object semantics. Inspired by this, we develop a zero-shot unsupervised approach for learning object segmentations. The model comprises two visual pathways: an appearance pathway that segments individual RGB images into coherent object regions, and a motion pathway that predicts the flow vector for each region between consecutive video frames. The two pathways jointly reconstruct a new representation called segment flow. This decoupled representation of appearance and motion is trained in a self-supervised manner to reconstruct one frame from another. When pretrained on an unlabeled video corpus, the model can be useful for a variety of applications, including 1) primary object segmentation from a single image in a zero-shot fashion; 2) moving object segmentation from a video with unsupervised test-time adaptation; 3) image semantic segmentation by supervised fine-tuning on a labeled image dataset. We demonstrate encouraging experimental results on all of these tasks using pretrained models.

[Direct Multi-view Multi-person 3D Pose Estimation](#)

- tao wang, Jianfeng Zhang, Yujun Cai, Shuicheng Yan, Jiashi Feng
- abstract@[open-review](#): We present Multi-view Pose transformer (MvP) for estimating multi-person 3D poses from multi-view images. Instead of estimating 3D joint locations from costly volumetric representation or reconstructing the per-person 3D pose from multiple detected 2D poses as in previous methods, MvP directly regresses the multi-person 3D poses in a clean and efficient way, without relying on intermediate tasks. Specifically, MvP represents skeleton joints as learnable query embeddings and let them progressively attend to and reason over the multi-view information from the input images to directly regress the actual 3D joint locations. To improve the accuracy of such a simple pipeline, MvP presents a hierarchical scheme to concisely represent query embeddings of multi-person skeleton joints and introduces an input-dependent query adaptation approach. Further, MvP designs a novel geometrically guided attention mechanism, called projective attention, to more precisely fuse the cross-view information for each joint. MvP also introduces a RayConv operation to integrate the view-dependent camera geometry into the feature representations for augmenting the projective attention. We show experimentally that our MvP model outperforms the state-of-the-art methods on several benchmarks while being much more efficient. Notably, it achieves 92.3% AP25 on the challenging Panoptic dataset, improving upon the previous best approach [35] by 9.8%. MvP is general and also extendable to recovering human mesh represented by the SMPL model, thus useful for modeling multi-person body shapes. Code and models are available at <https://github.com/sail-sg/mvp>.

[MST: Masked Self-Supervised Transformer for Visual Representation](#)

- Zhaowen Li, Zhiyang Chen, Fan Yang, Wei Li, YouSong Zhu, Chaoyang Zhao, Rui Deng, Liwei Wu, Rui Zhao, Ming Tang, Jinqiao Wang
- abstract@[open-review](#): Transformer has been widely used for self-supervised pre-training in Natural Language Processing (NLP) and achieved great success. However, it has not been fully explored in visual self-supervised learning. Meanwhile, previous methods only consider the high-level feature and learning representation from a global perspective, which may fail to transfer to the downstream dense prediction tasks focusing on local features. In this paper, we present a novel Masked Self-supervised Transformer approach named MST, which can explicitly capture the local context of an image while preserving the global semantic information. Specifically, inspired by the Masked Language Modeling (MLM) in NLP, we propose a masked token strategy based on the multi-head self-attention map, which dynamically masks some tokens of local patches without damaging the crucial structure for self-supervised learning. More importantly, the masked tokens together with the remaining tokens are further recovered by a global image decoder, which preserves the spatial information of the image and is more friendly to the downstream dense prediction tasks. The experiments on multiple datasets demonstrate the effectiveness and generality of the proposed method. For instance, MST achieves Top-1 accuracy of 76.9% with DeiT-S only using 300-epoch pre-training by linear evaluation, which outperforms supervised methods with the same epoch by 0.4% and its comparable variant DINO by 1.0%. For dense prediction tasks, MST also achieves 42.7% mAP on MS COCO object detection and 74.04% mIoU on Cityscapes segmentation only with 100-epoch pre-training.

[Exploiting Opponents Under Utility Constraints in Sequential Games](#)

- Martino Bernasconi-de-Luca, Federico Cacciamani, Simone Fioravanti, Nicola Gatti, Alberto Marchesi, Francesco Trovò
- abstract@[open-review](#): Recently, game-playing agents based on AI techniques have demonstrated super-human performance in several sequential games, such as chess, Go, and poker. Surprisingly, the multi-agent learning techniques that allowed to reach these achievements do not take into account the actual behavior of the human player, potentially leading to an impressive gap in performances. In this paper, we address the problem of designing artificial agents that learn how to effectively exploit unknown human opponents while playing repeatedly against them in an online fashion. We study the case in

which the agent's strategy during each repetition of the game is subject to constraints ensuring that the human's expected utility is within some lower and upper thresholds. Our framework encompasses several real-world problems, such as human engagement in repeated game playing and human education by means of serious games. As a first result, we formalize a set of linear inequalities encoding the conditions that the agent's strategy must satisfy at each iteration in order to do not violate the given bounds for the human's expected utility. Then, we use such formulation in an upper confidence bound algorithm, and we prove that the resulting procedure suffers from sublinear regret and guarantees that the constraints are satisfied with high probability at each iteration. Finally, we empirically evaluate the convergence of our algorithm on standard testbeds of sequential games.

[A Compositional Atlas of Tractable Circuit Operations for Probabilistic Inference](#)

- Antonio Vergari, YooJung Choi, Anji Liu, Stefano Teso, Guy Van den Broeck
- abstract@[open-review](#): Circuit representations are becoming the lingua franca to express and reason about tractable generative and discriminative models. In this paper, we show how complex inference scenarios for these models that commonly arise in machine learning---from computing the expectations of decision tree ensembles to information-theoretic divergences of sum-product networks---can be represented in terms of tractable modular operations over circuits. Specifically, we characterize the tractability of simple transformations---sums, products, quotients, powers, logarithms, and exponentials---in terms of sufficient structural constraints of the circuits they operate on, and present novel hardness results for the cases in which these properties are not satisfied. Building on these operations, we derive a unified framework for reasoning about tractable models that generalizes several results in the literature and opens up novel tractable inference scenarios.

[Demystifying and Generalizing BinaryConnect](#)

- Tim Dolkhorn, Yaoliang Yu, Eyyüb Sari, Mahdi Zolnouri, Vahid Partovi Nia
- abstract@[open-review](#): BinaryConnect (BC) and its many variations have become the de facto standard for neural network quantization. However, our understanding of the inner workings of BC is still quite limited. We attempt to close this gap in four different aspects: (a) we show that existing quantization algorithms, including post-training quantization, are surprisingly similar to each other; (b) we argue for proximal maps as a natural family of quantizers that is both easy to design and analyze; (c) we refine the observation that BC is a special case of dual averaging, which itself is a special case of the generalized conditional gradient algorithm; (d) consequently, we propose ProxConnect (PC) as a generalization of BC and we prove its convergence properties by exploiting the established connections. We conduct experiments on CIFAR-10 and ImageNet, and verify that PC achieves competitive performance.

[CARMS: Categorical-Antithetic-REINFORCE Multi-Sample Gradient Estimator](#)

- Alek Dimitriev, Mingyuan Zhou
- abstract@[open-review](#): Accurately backpropagating the gradient through categorical variables is a challenging task that arises in various domains, such as training discrete latent variable models. To this end, we propose CARMS, an unbiased estimator for categorical random variables based on multiple mutually negatively correlated (jointly antithetic) samples. CARMS combines REINFORCE with copula based sampling to avoid duplicate samples and reduce its variance, while keeping the estimator unbiased using importance sampling. It generalizes both the ARMS antithetic estimator for binary variables, which is CARMS for two categories, as well as LOORF/VarGrad, the leave-one-out REINFORCE estimator, which is CARMS with independent samples. We evaluate CARMS on several benchmark datasets on a generative modeling task, as well as a structured output prediction task, and find it to outperform competing methods including a strong self-control baseline. The code is publicly available.

[Learning to Learn Dense Gaussian Processes for Few-Shot Learning](#)

- Ze Wang, Zichen Miao, Xiantong Zhen, Qiang Qiu
- abstract@[open-review](#): Gaussian processes with deep neural networks demonstrate to be a strong learner for few-shot learning since they combine the strength of deep learning and kernels while being able to well capture uncertainty. However, it remains an open problem to leverage the shared knowledge provided by related tasks. In this paper, we propose to learn Gaussian processes with dense inducing variables by meta-learning for few-shot learning. In contrast to sparse Gaussian processes, we define a set of dense inducing variables to be of a much larger size than the support set in each task, which collects prior knowledge from experienced tasks. The dense inducing variables specify a shared Gaussian process prior over prediction functions of all tasks, which are learned in a variational inference framework and offer a strong inductive bias for learning new tasks. To achieve task-specific prediction functions, we propose to adapt the inducing variables to each task by efficient gradient descent. We conduct extensive experiments on common benchmark datasets for a variety of few-shot learning tasks. Our dense Gaussian processes present significant improvements over vanilla Gaussian processes and comparable or even better performance with state-of-the-art methods.

[Stochastic Solutions for Linear Inverse Problems using the Prior Implicit in a Denoiser](#)

- Zahra Kadkhodaie, Eero Simoncelli
- abstract@[open-review](#): Deep neural networks have provided state-of-the-art solutions for problems such as image denoising, which implicitly rely on a prior probability model of natural images. Two recent lines of work – Denoising Score Matching and Plug-and-Play – propose methodologies for drawing samples from this implicit prior and using it to solve inverse problems, respectively. Here, we develop a parsimonious and robust generalization of these ideas. We rely on a classic statistical result that shows the least-squares solution for removing additive Gaussian noise can be written directly in terms of the gradient of the log of the noisy signal density. We use this to derive a stochastic coarse-to-fine gradient ascent procedure for drawing high-probability samples from the implicit prior embedded within a CNN trained to perform blind denoising. A generalization of this algorithm to constrained sampling provides a method for using the implicit prior to solve any deterministic linear inverse problem, with no additional training, thus extending the power of supervised learning for denoising to a much broader set of problems. The algorithm relies on minimal assumptions and exhibits robust convergence over a wide range of parameter choices. To demonstrate the generality of our method, we use it to obtain state-of-the-art levels of unsupervised performance for deblurring, super-resolution, and compressive sensing.

[Towards Stable and Robust AdderNets](#)

- Minjing Dong, Yunhe Wang, Xinghao Chen, Chang Xu
- abstract@[open-review](#): Adder neural network (AdderNet) replaces the original convolutions with massive multiplications by cheap additions while achieving comparable performance thus yields a series of energy-efficient neural networks. Compared with convolutional neural networks (CNNs), the training of AdderNets is much more sophisticated including several techniques for adjusting gradient and batch normalization. In addition, variances of both weights and activations in resulting adder networks are very enormous which limits its performance and the potential for applying to other tasks. To enhance the stability and robustness of AdderNets, we first thoroughly analyze the variance estimation of weight parameters and output features of an arbitrary adder layer. Then, we develop a weight normalization scheme for adaptively optimizing the weight distribution of AdderNets during the training procedure, which can reduce the perturbation on running mean and variance in batch normalization layers. Meanwhile, the proposed weight normalization can also be utilized to enhance the adversarial robustness of resulting networks. Experiments conducted on several benchmarks demonstrate the superiority of the proposed approach for generating AdderNets with higher performance.

[Representing Long-Range Context for Graph Neural Networks with Global Attention](#)

- Zhanghao Wu, Paras Jain, Matthew Wright, Azalia Mirhoseini, Joseph E. Gonzalez, Ion Stoica
- abstract@[open-review](#): Graph neural networks are powerful architectures for structured datasets. However, current methods struggle to represent long-range dependencies. Scaling the depth or width of GNNs is insufficient to broaden receptive fields as larger GNNs encounter optimization instabilities such as vanishing gradients and representation oversmoothing, while pooling-based approaches have yet to become as universally useful as in computer vision. In this work, we propose the use of Transformer-based self-attention to learn long-range pairwise relationships, with a novel “readout” mechanism to obtain a global graph embedding. Inspired by recent computer vision results that find position-invariant attention performant in learning long-range relationships, our method, which we call GraphTrans, applies a permutation-invariant Transformer module after a standard GNN module. This simple architecture leads to state-of-the-art results on several graph classification tasks, outperforming methods that explicitly encode graph structure. Our results suggest that purely-learning-based approaches without graph structure may be suitable for learning high-level, long-range relationships on graphs. Code for GraphTrans is available at <https://github.com/ucbrise/graphtrans>.

[Beyond Bandit Feedback in Online Multiclass Classification](#)

- Dirk van der Hoeven, Federico Fusco, Nicolò Cesa-Bianchi
- abstract@[open-review](#): We study the problem of online multiclass classification in a setting where the learner's feedback is determined by an arbitrary directed graph. While including bandit feedback as a special case, feedback graphs allow a much richer set of applications, including filtering and label efficient classification. We introduce \textproc{Gappletron}, the first online multiclass algorithm that works with arbitrary feedback graphs. For this new algorithm, we prove surrogate regret bounds that hold, both in expectation and with high probability, for a large class of surrogate losses. Our bounds are of order $\$B\sqrt{\rho KT}$, where B is the diameter of the prediction space, K is the number of classes, T is the time horizon, and ρ is the domination number (a graph-theoretic parameter affecting the amount of exploration). In the full information case, we show that \textproc{Gappletron} achieves a constant surrogate regret of order B^2K . We also prove a general lower bound of order $\max\{B^2K, \sqrt{T}\}$ showing that our upper bounds are not significantly improvable. Experiments on synthetic data show that for various feedback graphs our algorithm is competitive against known baselines.

[Learning Student-Friendly Teacher Networks for Knowledge Distillation](#)

- Dae Young Park, Moon-Hyun Cha, Changwook Jeong, Daesin Kim, Bohyung Han
- abstract@[open-review](#): We propose a novel knowledge distillation approach to facilitate the transfer of dark knowledge from a teacher to a student. Contrary to most of the existing methods that rely on effective training of student models given pretrained teachers, we aim to learn the teacher models that are friendly to students and, consequently, more appropriate for knowledge transfer. In other words, at the time of optimizing a teacher model, the proposed algorithm learns the student branches jointly to obtain student-friendly representations. Since the main goal of our approach lies in training teacher models and the subsequent knowledge distillation procedure is straightforward, most of the existing knowledge distillation methods can adopt this technique to improve the performance of diverse student models in terms of accuracy and convergence speed. The proposed algorithm demonstrates outstanding accuracy in several well-known knowledge distillation techniques with various combinations of teacher and student models even in the case that their architectures are heterogeneous and there is no prior knowledge about student models at the time of training teacher networks

[Implicit Transformer Network for Screen Content Image Continuous Super-Resolution](#)

- Jingyu Yang, Sheng Shen, Huanjing Yue, Kun Li
- abstract@[open-review](#): Nowadays, there is an explosive growth of screen contents due to the wide application of screen sharing, remote cooperation, and online education. To match the limited terminal bandwidth, high-resolution (HR) screen contents may be downsampled and compressed. At the receiver side, the super-resolution (SR) of low-resolution (LR) screen content images (SCIs) is highly demanded by the HR display or by the users to zoom in for detail observation. However, image SR methods mostly designed for natural images do not generalize well for SCIs due to the very different image characteristics as well as the requirement of SCI browsing at arbitrary scales. To this end, we propose a novel Implicit Transformer Super-Resolution Network (ITSRN) for SCISR. For high-quality continuous SR at arbitrary ratios, pixel values at query coordinates are inferred from image features at key coordinates by the proposed implicit transformer and an implicit position encoding scheme is proposed to aggregate similar neighboring pixel values to the query one. We construct benchmark SCI1K and SCI1K-compression datasets with LR and HR SCI pairs. Extensive experiments show that the proposed ITsrn significantly outperforms several competitive continuous and discrete SR methods for both compressed and uncompressed SCIs.

[Channel Permutations for N:M Sparsity](#)

- Jeff Pool, Chong Yu
- abstract@[open-review](#): We introduce channel permutations as a method to maximize the accuracy of N:M sparse networks. N:M sparsity requires N out of M consecutive elements to be zero and has been shown to maintain accuracy for many models and tasks with a simple prune and fine-tune workflow. By permuting weight matrices along their channel dimension and adjusting the surrounding layers appropriately, we demonstrate accuracy recovery for even small, parameter-efficient networks, without affecting inference run-time. We also present both a quality metric to simplify judging permutations as well as efficient methods to search for high-quality permutations, including two optimizations to escape local minima. Finally, we share an ablation study to show the importance of each part of our search algorithm, experimental results showing correlation between our quality metric and final network accuracy, improved sparse network accuracy using our techniques with insignificant overhead to training time, and the transformation of unstructured to structured sparse workloads. Code to use these techniques when generating a 2:4 sparse network is available at <https://github.com/NVIDIA/apex/tree/master/apex/contrib/sparsity>.

[Curriculum Learning for Vision-and-Language Navigation](#)

- Jiwen Zhang, Zhongyu Wei, Jianqing Fan, Jiajie Peng
- abstract@[open-review](#): Vision-and-Language Navigation (VLN) is a task where an agent navigates in an embodied indoor environment under human instructions. Previous works ignore the distribution of sample difficulty and we argue that this potentially degrades their agent performance. To tackle this issue, we propose a novel curriculum-based training paradigm for VLN tasks that can balance human prior knowledge and agent learning progress about training samples. We develop the principle of curriculum design and re-arrange the benchmark Room-to-Room (R2R) dataset to make it suitable for curriculum training. Experiments show that our method is model-agnostic and can significantly improve the performance, the generalizability, and the training efficiency of current state-of-the-art navigation agents without increasing model complexity.

[Better Algorithms for Individually Fair \$k\$ -Clustering](#)

- Maryam Negahbani, Deeparnab Chakrabarty
- abstract@[open-review](#): We study data clustering problems with ℓ_p -norm objectives (e.g. \textsc{\$k\$-Median} and \textsc{\$k\$-Means}) in the context of individual fairness. The dataset consists of n points, and we want to find k centers such that (a) the objective is minimized, while (b) respecting the individual fairness constraint that every point v has a center within a distance at most $r(v)$, where $r(v)$ is v 's distance to its (n/k) th nearest point. Jung, Kannan, and Lutz [FORC 2020] introduced this concept and designed a clustering algorithm with provable (approximate) fairness and objective guarantees for the ℓ_∞ or \textsc{\$k\$-Center} objective. Mahabadi and Vakilian [ICML 2020] revisited this problem to give a local-search algorithm for all ℓ_p -norms. Empirically, their algorithms outperform Jung et. al.'s by a large margin in terms of cost (for \textsc{\$k\$-Median} and \textsc{\$k\$-Means}), but they incur a reasonable loss in fairness. In this paper, our main contribution is to use Linear

Programming (LP) techniques to obtain better algorithms for this problem, both in theory and in practice. We prove that by modifying known LP rounding techniques, one gets a worst-case guarantee on the objective which is much better than in MV20, and empirically, this objective is extremely close to the optimal. Furthermore, our theoretical fairness guarantees are comparable with MV20 in theory, and empirically, we obtain noticeably fairer solutions. Although solving the LP $\{\text{em exactly}\}$ might be prohibitive, we demonstrate that in practice, a simple sparsification technique drastically improves the run-time of our algorithm.

[Video Instance Segmentation using Inter-Frame Communication Transformers](#)

- Sukjun Hwang, Miran Heo, Seoung Wug Oh, Seon Joo Kim
- abstract@[open-review](#): We propose a novel end-to-end solution for video instance segmentation (VIS) based on transformers. Recently, the per-clip pipeline shows superior performance over per-frame methods leveraging richer information from multiple frames. However, previous per-clip models require heavy computation and memory usage to achieve frame-to-frame communications, limiting practicality. In this work, we propose Inter-frame Communication Transformers (IFC), which significantly reduces the overhead for information-passing between frames by efficiently encoding the context within the input clip. Specifically, we propose to utilize concise memory tokens as a means of conveying information as well as summarizing each frame scene. The features of each frame are enriched and correlated with other frames through exchange of information between the precisely encoded memory tokens. We validate our method on the latest benchmark sets and achieved state-of-the-art performance (AP 42.6 on YouTube-VIS 2019 val set using the offline inference) while having a considerably fast runtime (89.4 FPS). Our method can also be applied to near-online inference for processing a video in real-time with only a small delay. The code is available at <https://github.com/sukjunhwang/IFC>

[Progressive Coordinate Transforms for Monocular 3D Object Detection](#)

- Li Wang, Li Zhang, Yi Zhu, Zhi Zhang, Tong He, Mu Li, Xiangyang Xue
- abstract@[open-review](#): Recognizing and localizing objects in the 3D space is a crucial ability for an AI agent to perceive its surrounding environment. While significant progress has been achieved with expensive LiDAR point clouds, it poses a great challenge for 3D object detection given only a monocular image. While there exist different alternatives for tackling this problem, it is found that they are either equipped with heavy networks to fuse RGB and depth information or empirically ineffective to process millions of pseudo-LiDAR points. With in-depth examination, we realize that these limitations are rooted in inaccurate object localization. In this paper, we propose a novel and lightweight approach, dubbed {\\em Progressive Coordinate Transforms} (PCT) to facilitate learning coordinate representations. Specifically, a localization boosting mechanism with confidence-aware loss is introduced to progressively refine the localization prediction. In addition, semantic image representation is also exploited to compensate for the usage of patch proposals. Despite being lightweight and simple, our strategy allows us to establish a new state-of-the-art among the monocular 3D detectors on the competitive KITTI benchmark. At the same time, our proposed PCT shows great generalization to most coordinate-based 3D detection frameworks.

[Structured Reordering for Modeling Latent Alignments in Sequence Transduction](#)

- bailin wang, Mirella Lapata, Ivan Titov
- abstract@[open-review](#): Despite success in many domains, neural models struggle in settings where train and test examples are drawn from different distributions. In particular, in contrast to humans, conventional sequence-to-sequence (seq2seq) models fail to generalize systematically, i.e., interpret sentences representing novel combinations of concepts (e.g., text segments) seen in training. Traditional grammar formalisms excel in such settings by implicitly encoding alignments between input and output segments, but are hard to scale and maintain. Instead of engineering a grammar, we directly model segment-to-segment alignments as discrete structured latent variables within a neural seq2seq model. To efficiently explore the large space of alignments, we introduce a reorder-first align-later framework whose central component is a neural reordering module producing separable permutations. We present an efficient dynamic programming algorithm performing exact marginal inference of separable permutations, and, thus, enabling end-to-end differentiable training of our model. The resulting seq2seq model exhibits better systematic generalization than standard models on synthetic problems and NLP tasks (i.e., semantic parsing and machine translation).

[A universal probabilistic spike count model reveals ongoing modulation of neural variability](#)

- David Liu, Mate Lengyel
- abstract@[open-review](#): Neural responses are variable: even under identical experimental conditions, single neuron and population responses typically differ from trial to trial and across time. Recent work has demonstrated that this variability has predictable structure, can be modulated by sensory input and behaviour, and bears critical signatures of the underlying network dynamics and computations. However, current methods for characterising neural variability are primarily geared towards sensory coding in the laboratory: they require trials with repeatable experimental stimuli and behavioural covariates. In addition, they make strong assumptions about the parametric form of variability, rely on assumption-free but data-inefficient histogram-based approaches, or are altogether ill-suited for capturing variability modulation by covariates. Here we present a universal probabilistic spike count model that eliminates these shortcomings. Our method builds on sparse Gaussian processes and can model arbitrary spike count distributions (SCDs) with flexible dependence on observed as well as latent covariates, using scalable variational inference to jointly infer the covariate-to-SCD mappings and latent trajectories in a data efficient way. Without requiring repeatable trials, it can flexibly capture covariate-dependent joint SCDs, and provide interpretable latent causes underlying the statistical dependencies between neurons. We apply the model to recordings from a canonical non-sensory neural population: head direction cells in the mouse. We find that variability in these cells defies a simple parametric relationship with mean spike count as assumed in standard models, its modulation by external covariates can be comparably strong to that of the mean firing rate, and slow low-dimensional latent factors explain away neural correlations. Our approach paves the way to understanding the mechanisms and computations underlying neural variability under naturalistic conditions, beyond the realm of sensory coding with repeatable stimuli.

[Bellman Eluder Dimension: New Rich Classes of RL Problems, and Sample-Efficient Algorithms](#)

- Chi Jin, Qinghua Liu, Sobhan Miryoosefi
- abstract@[open-review](#): Finding the minimal structural assumptions that empower sample-efficient learning is one of the most important research directions in Reinforcement Learning (RL). This paper advances our understanding of this fundamental question by introducing a new complexity measure—Bellman Eluder (BE) dimension. We show that the family of RL problems of low BE dimension is remarkably rich, which subsumes a vast majority of existing tractable RL problems including but not limited to tabular MDPs, linear MDPs, reactive POMDPs, low Bellman rank problems as well as low Eluder dimension problems. This paper further designs a new optimization-based algorithm—GOLF, and reanalyzes a hypothesis elimination-based algorithm—OLIVE (proposed in Jiang et al. (2017)). We prove that both algorithms learn the near-optimal policies of low BE dimension problems in a number of samples that is polynomial in all relevant parameters, but independent of the size of state-action space. Our regret and sample complexity results match or improve the best existing results for several well-known subclasses of low BE dimension problems.

[Detecting Anomalous Event Sequences with Temporal Point Processes](#)

- Oleksandr Shchur, Ali Caner Turkmen, Tim Januschowski, Jan Gasthaus, Stephan Günnemann
- abstract@[open-review](#): Automatically detecting anomalies in event data can provide substantial value in domains such as healthcare, DevOps, and information security. In this paper, we frame the problem of detecting anomalous continuous-time event sequences as out-of-distribution (OOD) detection for temporal point processes (TPPs). First, we show how this problem can be approached using goodness-of-fit (GoF) tests. We then demonstrate the limitations of popular GoF statistics for TPPs and propose a new test that addresses these shortcomings. The proposed method can be combined with

various TPP models, such as neural TPPs, and is easy to implement. In our experiments, we show that the proposed statistic excels at both traditional GoF testing, as well as at detecting anomalies in simulated and real-world data.

[HNPE: Leveraging Global Parameters for Neural Posterior Estimation](#)

- Pedro Rodrigues, Thomas Moreau, Gilles Louppe, Alexandre Gramfort
- abstract@[open-review](#): Inferring the parameters of a stochastic model based on experimental observations is central to the scientific method. A particularly challenging setting is when the model is strongly indeterminate, i.e. when distinct sets of parameters yield identical observations. This arises in many practical situations, such as when inferring the distance and power of a radio source (is the source close and weak or far and strong?) or when estimating the amplifier gain and underlying brain activity of an electrophysiological experiment. In this work, we present hierarchical neural posterior estimation (HNPE), a novel method for cracking such indeterminacy by exploiting additional information conveyed by an auxiliary set of observations sharing global parameters. Our method extends recent developments in simulation-based inference (SBI) based on normalizing flows to Bayesian hierarchical models. We validate quantitatively our proposal on a motivating example amenable to analytical solutions and then apply it to invert a well known non-linear model from computational neuroscience, using both simulated and real EEG data.

[Alignment Attention by Matching Key and Query Distributions](#)

- Shujian Zhang, Xinjie Fan, Huangjie Zheng, Korawat Tanwisuth, Mingyuan Zhou
- abstract@[open-review](#): The neural attention mechanism has been incorporated into deep neural networks to achieve state-of-the-art performance in various domains. Most such models use multi-head self-attention which is appealing for the ability to attend to information from different perspectives. This paper introduces alignment attention that explicitly encourages self-attention to match the distributions of the key and query within each head. The resulting alignment attention networks can be optimized as an unsupervised regularization in the existing attention framework. It is simple to convert any models with self-attention, including pre-trained ones, to the proposed alignment attention. On a variety of language understanding tasks, we show the effectiveness of our method in accuracy, uncertainty estimation, generalization across domains, and robustness to adversarial attacks. We further demonstrate the general applicability of our approach on graph attention and visual question answering, showing the great potential of incorporating our alignment method into various attention-related tasks.

[Settling the Variance of Multi-Agent Policy Gradients](#)

- Jakub Grudzien Kuba, Muning Wen, Linghui Meng, shangding gu, Haifeng Zhang, David Mguni, Jun Wang, Yaodong Yang
- abstract@[open-review](#): Policy gradient (PG) methods are popular reinforcement learning (RL) methods where a baseline is often applied to reduce the variance of gradient estimates. In multi-agent RL (MARL), although the PG theorem can be naturally extended, the effectiveness of multi-agent PG (MAPG) methods degrades as the variance of gradient estimates increases rapidly with the number of agents. In this paper, we offer a rigorous analysis of MAPG methods by, firstly, quantifying the contributions of the number of agents and agents' explorations to the variance of MAPG estimators. Based on this analysis, we derive the optimal baseline (OB) that achieves the minimal variance. In comparison to the OB, we measure the excess variance of existing MARL algorithms such as vanilla MAPG and COMA. Considering using deep neural networks, we also propose a surrogate version of OB, which can be seamlessly plugged into any existing PG methods in MARL. On benchmarks of Multi-Agent MuJoCo and StarCraft challenges, our OB technique effectively stabilises training and improves the performance of multi-agent PPO and COMA algorithms by a significant margin. Code is released at \url{https://github.com/morning9393/Optimal-Baseline-for-Multi-agent-Policy-Gradients}.

[For high-dimensional hierarchical models, consider exchangeability of effects across covariates instead of across datasets](#)

- Brian Trippe, Hilary Finucane, Tamara Broderick
- abstract@[open-review](#): Hierarchical Bayesian methods enable information sharing across regression problems on multiple groups of data. While standard practice is to model regression parameters (effects) as (1) exchangeable across the groups and (2) correlated to differing degrees across covariates, we show that this approach exhibits poor statistical performance when the number of covariates exceeds the number of groups. For instance, in statistical genetics, we might regress dozens of traits (defining groups) for thousands of individuals (responses) on up to millions of genetic variants (covariates). When an analyst has more covariates than groups, we argue that it is often preferable to instead model effects as (1) exchangeable across covariates and (2) correlated to differing degrees across groups. To this end, we propose a hierarchical model expressing our alternative perspective. We devise an empirical Bayes estimator for learning the degree of correlation between groups. We develop theory that demonstrates that our method outperforms the classic approach when the number of covariates dominates the number of groups, and corroborate this result empirically on several high-dimensional multiple regression and classification problems.

[Efficient Algorithms for Learning Depth-2 Neural Networks with General ReLU Activations](#)

- Pranjal Awasthi, Alex Tang, Aravindan Vijayaraghavan
- abstract@[open-review](#): We present polynomial time and sample efficient algorithms for learning an unknown depth-2 feedforward neural network with general ReLU activations, under mild non-degeneracy assumptions. In particular, we consider learning an unknown network of the form $f(x) = \{a\}^{\mathsf{T}} \{W\} \sigma(\{T\}x + b)$, where x is drawn from the Gaussian distribution, and $\sigma(t) = \max(t, 0)$ is the ReLU activation. Prior works for learning networks with ReLU activations assume that the bias (b) is zero. In order to deal with the presence of the bias terms, our proposed algorithm consists of robustly decomposing multiple higher order tensors arising from the Hermite expansion of the function $f(x)$. Using these ideas we also establish identifiability of the network parameters under very mild assumptions.

[Controllable and Compositional Generation with Latent-Space Energy-Based Models](#)

- Weili Nie, Arash Vahdat, Anima Anandkumar
- abstract@[open-review](#): Controllable generation is one of the key requirements for successful adoption of deep generative models in real-world applications, but it still remains as a great challenge. In particular, the compositional ability to generate novel concept combinations is out of reach for most current models. In this work, we use energy-based models (EBMs) to handle compositional generation over a set of attributes. To make them scalable to high-resolution image generation, we introduce an EBM in the latent space of a pre-trained generative model such as StyleGAN. We propose a novel EBM formulation representing the joint distribution of data and attributes together, and we show how sampling from it is formulated as solving an ordinary differential equation (ODE). Given a pre-trained generator, all we need for controllable generation is to train an attribute classifier. Sampling with ODEs is done efficiently in the latent space and is robust to hyperparameters. Thus, our method is simple, fast to train, and efficient to sample. Experimental results show that our method outperforms the state-of-the-art in both conditional sampling and sequential editing. In compositional generation, our method excels at zero-shot generation of unseen attribute combinations. Also, by composing energy functions with logical operators, this work is the first to achieve such compositionality in generating photo-realistic images of resolution 1024x1024.

[Reverse-Complement Equivariant Networks for DNA Sequences](#)

- Vincent Mallet, Jean-Philippe Vert

- abstract@[open-review](#): As DNA sequencing technologies keep improving in scale and cost, there is a growing need to develop machine learning models to analyze DNA sequences, e.g., to decipher regulatory signals from DNA fragments bound by a particular protein of interest. As a double helix made of two complementary strands, a DNA fragment can be sequenced as two equivalent, so-called reverse complement (RC) sequences of nucleotides. To take into account this inherent symmetry of the data in machine learning models can facilitate learning. In this sense, several authors have recently proposed particular RC-equivariant convolutional neural networks (CNNs). However, it remains unknown whether other RC-equivariant architecture exist, which could potentially increase the set of basic models adapted to DNA sequences for practitioners. Here, we close this gap by characterizing the set of all linear RC-equivariant layers, and show in particular that new architectures exist beyond the ones already explored. We further discuss RC-equivariant pointwise nonlinearities adapted to different architectures, as well as RC-equivariant embeddings of k -mers as an alternative to one-hot encoding of nucleotides. We show experimentally that the new architectures can outperform existing ones.

[Provably Efficient Reinforcement Learning with Linear Function Approximation under Adaptivity Constraints](#)

- Tianhao Wang, Dongruo Zhou, Quanquan Gu
- abstract@[open-review](#): We study reinforcement learning (RL) with linear function approximation under the adaptivity constraint. We consider two popular limited adaptivity models: the batch learning model and the rare policy switch model, and propose two efficient online RL algorithms for episodic linear Markov decision processes, where the transition probability and the reward function can be represented as a linear function of some known feature mapping. In specific, for the batch learning model, our proposed LSVI-UCB-Batch algorithm achieves an $\tilde{O}(\sqrt{d^3H^3T} + dHT/B)$ regret, where d is the dimension of the feature mapping, H is the episode length, T is the number of interactions and B is the number of batches. Our result suggests that it suffices to use only $\sqrt{T/dH}$ batches to obtain $\tilde{O}(\sqrt{d^3H^3T})$ regret. For the rare policy switch model, our proposed LSVI-UCB-RareSwitch algorithm enjoys an $\tilde{O}(\sqrt{d^3H^3T[1+T/(dH)]^dH/B})$ regret, which implies that $dH \log T$ policy switches suffice to obtain the $\tilde{O}(\sqrt{d^3H^3T})$ regret. Our algorithms achieve the same regret as the LSVI-UCB algorithm [\[jin2020provably\]](#), yet with a substantially smaller amount of adaptivity. We also establish a lower bound for the batch learning model, which suggests that the dependency on B in our regret bound is tight.

[Nonsmooth Implicit Differentiation for Machine-Learning and Optimization](#)

- Jérôme Bolte, Tam Le, Edouard Pauwels, Tony Silvetti-Falls
- abstract@[open-review](#): In view of training increasingly complex learning architectures, we establish a nonsmooth implicit function theorem with an operational calculus. Our result applies to most practical problems (i.e., definable problems) provided that a nonsmooth form of the classical invertibility condition is fulfilled. This approach allows for formal subdifferentiation: for instance, replacing derivatives by Clarke Jacobians in the usual differentiation formulas is fully justified for a wide class of nonsmooth problems. Moreover this calculus is entirely compatible with algorithmic differentiation (e.g., backpropagation). We provide several applications such as training deep equilibrium networks, training neural nets with conic optimization layers, or hyperparameter-tuning for nonsmooth Lasso-type models. To show the sharpness of our assumptions, we present numerical experiments showcasing the extremely pathological gradient dynamics one can encounter when applying implicit algorithmic differentiation without any hypothesis.

[Heuristic-Guided Reinforcement Learning](#)

- Ching-An Cheng, Andrey Kolobov, Adith Swaminathan
- abstract@[open-review](#): We provide a framework to accelerate reinforcement learning (RL) algorithms by heuristics that are constructed by domain knowledge or offline data. Tabula rasa RL algorithms require environment interactions or computation that scales with the horizon of the sequential decision-making task. Using our framework, we show how heuristic-guided RL induces a much shorter horizon sub-problem that provably solves the original task. Our framework can be viewed as a horizon-based regularization for controlling bias and variance in RL under a finite interaction budget. In theory, we characterize the properties of a good heuristic and the resulting impact on RL acceleration. In particular, we introduce the novel concept of an improvable heuristic that can allow any RL agent to conservatively extrapolate beyond its prior knowledge. In practice, we instantiate our framework to accelerate several state-of-the-art algorithms in simulated robotic control tasks and procedurally generated games. Our framework complements the rich literature on warm-starting RL using expert demonstrations or exploratory data-sets, and creates a unified channel to inject prior knowledge into RL.

[Statistical Undecidability in Linear, Non-Gaussian Causal Models in the Presence of Latent Confounders](#)

- Konstantin Genin
- abstract@[open-review](#): If causal relationships are linear and acyclic and noise terms are independent and Gaussian, causal orientation is not identified from observational data --- even if faithfulness is satisfied (Spirtes et al., 2002). Shimizu et al. (2006) showed that acyclic, linear, {\bf non}-Gaussian (LiNGAM) causal models {\em are} identified from observational data, so long as no latent confounders are present. That holds even when faithfulness fails. Genin and Mayo-Wilson (2020) refine that result: not only are causal relationships identified, but causal orientation is {\em statistically decidable}. That means that for every $\epsilon > 0$ there is a method that converges in probability to the correct orientation and, at every sample size, outputs an incorrect orientation with probability less than ϵ . These results naturally raise questions about what happens in the presence of latent confounders. Hoyer et al. (2008) and Salehkaleybar et al. (2020) show that, although the causal model is not uniquely identified, causal orientation among observed variables is identified in the presence of latent confounders, so long as faithfulness is satisfied. This paper refines these results: although it is possible to converge to the right orientation in the limit, causal orientation is no longer statistically decidable---it is not possible to converge to the correct orientation with finite-sample bounds on the probability of orientation errors, even if faithfulness is satisfied. However, that limiting result suggests several adjustments to the LiNGAM model that may recover decidability.

[A novel notion of barycenter for probability distributions based on optimal weak mass transport](#)

- Elsa Cazelles, Felipe Tobar, Joaquin Fontbona
- abstract@[open-review](#): We introduce weak barycenters of a family of probability distributions, based on the recently developed notion of optimal weak transport of mass by Gozlan et al. (2017) and Backhoff-Veraguas et al. (2020). We provide a theoretical analysis of this object and discuss its interpretation in the light of convex ordering between probability measures. In particular, we show that, rather than averaging the input distributions in a geometric way (as the Wasserstein barycenter based on classic optimal transport does) weak barycenters extract common geometric information shared by all the input distributions, encoded as a latent random variable that underlies all of them. We also provide an iterative algorithm to compute a weak barycenter for a finite family of input distributions, and a stochastic algorithm that computes them for arbitrary populations of laws. The latter approach is particularly well suited for the streaming setting, i.e., when distributions are observed sequentially. The notion of weak barycenter and our approaches to compute it are illustrated on synthetic examples, validated on 2D real-world data and compared to standard Wasserstein barycenters.

[Temporal-attentive Covariance Pooling Networks for Video Recognition](#)

- Zilin Gao, Qilong Wang, Bingbing Zhang, Qinghua Hu, Peihua Li
- abstract@[open-review](#): For video recognition task, a global representation summarizing the whole contents of the video snippets plays an important role for the final performance. However, existing video architectures usually generate it by using a simple, global average pooling (GAP) method, which has limited ability to capture complex dynamics of videos. For image recognition task, there exist evidences showing that covariance pooling has stronger representation ability than GAP. Unfortunately, such plain covariance pooling used in image recognition is an orderless representative, which cannot

model spatio-temporal structure inherent in videos. Therefore, this paper proposes a Temporal-attentive Covariance Pooling (TCP), inserted at the end of deep architectures, to produce powerful video representations. Specifically, our TCP first develops a temporal attention module to adaptively calibrate spatio-temporal features for the succeeding covariance pooling, approximatively producing attentive covariance representations. Then, a temporal covariance pooling performs temporal pooling of the attentive covariance representations to characterize both intra-frame correlations and inter-frame cross-correlations of the calibrated features. As such, the proposed TCP can capture complex temporal dynamics. Finally, a fast matrix power normalization is introduced to exploit geometry of covariance representations. Note that our TCP is model-agnostic and can be flexibly integrated into any video architectures, resulting in TCPNet for effective video recognition. The extensive experiments on six benchmarks (e.g., Kinetics, Something-Something V1 and Charades) using various video architectures show our TCPNet is clearly superior to its counterparts, while having strong generalization ability. The source code is publicly available.

[Revisiting Smoothed Online Learning](#)

- Lijun Zhang, Wei Jiang, Shiyin Lu, Tianbao Yang
- abstract@[open-review](#): In this paper, we revisit the problem of smoothed online learning, in which the online learner suffers both a hitting cost and a switching cost, and target two performance metrics: competitive ratio and dynamic regret with switching cost. To bound the competitive ratio, we assume the hitting cost is known to the learner in each round, and investigate the simple idea of balancing the two costs by an optimization problem. Surprisingly, we find that minimizing the hitting cost alone is $\max(1, \frac{2}{\alpha})$ -competitive for α -polyhedral functions and $1 + \frac{4}{\lambda}$ -competitive for λ -quadratic growth functions, both of which improve state-of-the-art results significantly. Moreover, when the hitting cost is both convex and λ -quadratic growth, we reduce the competitive ratio to $1 + \frac{2}{\sqrt{\lambda}}$ by minimizing the weighted sum of the hitting cost and the switching cost. To bound the dynamic regret with switching cost, we follow the standard setting of online convex optimization, in which the hitting cost is convex but hidden from the learner before making predictions. We modify Ader, an existing algorithm designed for dynamic regret, slightly to take into account the switching cost when measuring the performance. The proposed algorithm, named as Smoothed Ader, attains an optimal $O(\sqrt{T(1+P_T)})$ bound for dynamic regret with switching cost, where P_T is the path-length of the comparator sequence. Furthermore, if the hitting cost is accessible in the beginning of each round, we obtain a similar guarantee without the bounded gradient condition, and establish an $\Omega(\sqrt{T(1+P_T)})$ lower bound to confirm the optimality.

[Marginalised Gaussian Processes with Nested Sampling](#)

- Fergus Simpson, Vidhi Lalchand, Carl Edward Rasmussen
- abstract@[open-review](#): Gaussian Process models are a rich distribution over functions with inductive biases controlled by a kernel function. Learning occurs through optimisation of the kernel hyperparameters using the marginal likelihood as the objective. This work proposes nested sampling as a means of marginalising kernel hyperparameters, because it is a technique that is well-suited to exploring complex, multi-modal distributions. We benchmark against Hamiltonian Monte Carlo on time-series and two-dimensional regression tasks, finding that a principled approach to quantifying hyperparameter uncertainty substantially improves the quality of prediction intervals.

[Provable Benefits of Actor-Critic Methods for Offline Reinforcement Learning](#)

- Andrea Zanette, Martin J Wainwright, Emma Brunskill
- abstract@[open-review](#): Actor-critic methods are widely used in offline reinforcement learning practice, but are not so well-understood theoretically. We propose a new offline actor-critic algorithm that naturally incorporates the pessimism principle, leading to several key advantages compared to the state of the art. The algorithm can operate when the Bellman evaluation operator is closed with respect to the action value function of the actor's policies; this is a more general setting than the low-rank MDP model. Despite the added generality, the procedure is computationally tractable as it involves the solution of a sequence of second-order programs. We prove an upper bound on the suboptimality gap of the policy returned by the procedure that depends on the data coverage of any arbitrary, possibly data dependent comparator policy. The achievable guarantee is complemented with a minimax lower bound that is matching up to logarithmic factors.

[Bayesian Bellman Operators](#)

- Mattie Fellows, Kristian Hartikainen, Shimon Whiteson
- abstract@[open-review](#): We introduce a novel perspective on Bayesian reinforcement learning (RL); whereas existing approaches infer a posterior over the transition distribution or Q-function, we characterise the uncertainty in the Bellman operator. Our Bayesian Bellman operator (BBO) framework is motivated by the insight that when bootstrapping is introduced, model-free approaches actually infer a posterior over Bellman operators, not value functions. In this paper, we use BBO to provide a rigorous theoretical analysis of model-free Bayesian RL to better understand its relationship to established frequentist RL methodologies. We prove that Bayesian solutions are consistent with frequentist RL solutions, even when approximate inference is used, and derive conditions for which convergence properties hold. Empirically, we demonstrate that algorithms derived from the BBO framework have sophisticated deep exploration properties that enable them to solve continuous control tasks at which state-of-the-art regularised actor-critic algorithms fail catastrophically.

[Uncertainty Calibration for Ensemble-Based Debiasing Methods](#)

- Ruibin Xiong, Yimeng Chen, Liang Pang, Xueqi Cheng, Zhi-Ming Ma, Yanyan Lan
- abstract@[open-review](#): Ensemble-based debiasing methods have been shown effective in mitigating the reliance of classifiers on specific dataset bias, by exploiting the output of a bias-only model to adjust the learning target. In this paper, we focus on the bias-only model in these ensemble-based methods, which plays an important role but has not gained much attention in the existing literature. Theoretically, we prove that the debiasing performance can be damaged by inaccurate uncertainty estimations of the bias-only model. Empirically, we show that existing bias-only models fall short in producing accurate uncertainty estimations. Motivated by these findings, we propose to conduct calibration on the bias-only model, thus achieving a three-stage ensemble-based debiasing framework, including bias modeling, model calibrating, and debiasing. Experimental results on NLI and fact verification tasks show that our proposed three-stage debiasing framework consistently outperforms the traditional two-stage one in out-of-distribution accuracy.

[Provably Faster Algorithms for Bilevel Optimization](#)

- Junjie Yang, Kaiyi Ji, Yingbin Liang
- abstract@[open-review](#): Bilevel optimization has been widely applied in many important machine learning applications such as hyperparameter optimization and meta-learning. Recently, several momentum-based algorithms have been proposed to solve bilevel optimization problems faster. However, those momentum-based algorithms do not achieve provably better computational complexity than $\tilde{O}(\epsilon^{-2})$ of the SGD-based algorithm. In this paper, we propose two new algorithms for bilevel optimization, where the first algorithm adopts momentum-based recursive iterations, and the second algorithm adopts recursive gradient estimations in nested loops to decrease the variance. We show that both algorithms achieve the complexity of $\tilde{O}(\epsilon^{-1.5})$, which outperforms all existing algorithms by the order of magnitude. Our experiments validate our theoretical results and demonstrate the superior empirical performance of our algorithms in hyperparameter applications.

[Neo-GNNs: Neighborhood Overlap-aware Graph Neural Networks for Link Prediction](#)

- Seongjun Yun, Seoyoon Kim, Junhyun Lee, Jaewoo Kang, Hyunwoo J. Kim
- abstract@[open-review](#): Graph Neural Networks (GNNs) have been widely applied to various fields for learning over graph-structured data. They have shown significant improvements over traditional heuristic methods in various tasks such as node classification and graph classification. However, since GNNs heavily rely on smoothed node features rather than graph structure, they often show poor performance than simple heuristic methods in link prediction where the structural information, e.g., overlapped neighborhoods, degrees, and shortest paths, is crucial. To address this limitation, we propose Neighborhood Overlap-aware Graph Neural Networks (Neo-GNNs) that learn useful structural features from an adjacency matrix and estimate overlapped neighborhoods for link prediction. Our Neo-GNNs generalize neighborhood overlap-based heuristic methods and handle overlapped multi-hop neighborhoods. Our extensive experiments on Open Graph Benchmark datasets (OGB) demonstrate that Neo-GNNs consistently achieve state-of-the-art performance in link prediction.

[Self-Supervised Multi-Object Tracking with Cross-input Consistency](#)

- Favyen Bastani, Songtao He, Samuel Madden
- abstract@[open-review](#): In this paper, we propose a self-supervised learning procedure for training a robust multi-object tracking (MOT) model given only unlabeled video. While several self-supervisory learning signals have been proposed in prior work on single-object tracking, such as color propagation and cycle-consistency, these signals are not effective for training RNN models, which are needed to achieve accurate MOT: they yield degenerate models that, for instance, always match new detections to tracks with the closest initial detections. We propose a novel self-supervisory signal that we call cross-input consistency: we construct two distinct inputs for the same sequence of video, by hiding different information about the sequence in each input. We then compute tracks in that sequence by applying an RNN model independently on each input, and train the model to produce consistent tracks across the two inputs. We evaluate our unsupervised method on MOT17 and KITTI --- remarkably, we find that, despite training only on unlabeled video, our unsupervised approach outperforms four supervised methods published in the last 1--2 years, including Tracktor++, FAMNet, GSM, and mmMOT.

[Tree in Tree: from Decision Trees to Decision Graphs](#)

- Bingzhao Zhu, Mahsa Shoaran
- abstract@[open-review](#): Decision trees have been widely used as classifiers in many machine learning applications thanks to their lightweight and interpretable decision process. This paper introduces Tree in Tree decision graph (TnT), a framework that extends the conventional decision tree to a more generic and powerful directed acyclic graph. TnT constructs decision graphs by recursively growing decision trees inside the internal or leaf nodes instead of greedy training. The time complexity of TnT is linear to the number of nodes in the graph, therefore it can construct decision graphs on large datasets. Compared to decision trees, we show that TnT achieves better classification performance with reduced model size, both as a stand-alone classifier and as a base-estimator in bagging/AdaBoost ensembles. Our proposed model is a novel, more efficient and accurate alternative to the widely-used decision trees.

[Test-time Collective Prediction](#)

- Celestine Mendler-Dünner, Wenshuo Guo, Stephen Bates, Michael Jordan
- abstract@[open-review](#): An increasingly common setting in machine learning involves multiple parties, each with their own data, who want to jointly make predictions on future test points. Agents wish to benefit from the collective expertise of the full set of agents to make better predictions than they would individually, but may not be willing to release labeled data or model parameters. In this work, we explore a decentralized mechanism to make collective predictions at test time, that is inspired by the literature in social science on human consensus-making. Building on a query model to facilitate information exchange among agents, our approach leverages each agent's pre-trained model without relying on external validation, model retraining, or data pooling. A theoretical analysis shows that our approach recovers inverse mean-squared-error (MSE) weighting in the large-sample limit which is known to be the optimal way to combine independent, unbiased estimators. Empirically, we demonstrate that our scheme effectively combines models with differing quality across the input space: the proposed consensus prediction achieves significant gains over classical model averaging, and even outperforms weighted averaging schemes that have access to additional validation data. Finally, we propose a decentralized Jackknife procedure as a tool to evaluate the sensitivity of the collective predictions with respect to a single agent's opinion.

[A Continuous Mapping For Augmentation Design](#)

- Keyu Tian, Chen Lin, Ser Nam Lim, Wanli Ouyang, Puneet Dokania, Philip Torr
- abstract@[open-review](#): Automated data augmentation (ADA) techniques have played an important role in boosting the performance of deep models. Such techniques mostly aim to optimize a parameterized distribution over a discrete augmentation space. Thus, are restricted by the discretization of the search space which normally is handcrafted. To overcome the limitations, we take the first step to constructing a continuous mapping from \mathbb{R}^d to image transformations (an augmentation space). Using this mapping, we take a novel approach where 1) we pose the ADA as a continuous optimization problem over the parameters of the augmentation distribution; and 2) use Stochastic Gradient Langevin Dynamics to learn and sample augmentations. This allows us to potentially explore the space of infinitely many possible augmentations, which otherwise was not possible due to the discretization of the space. This view of ADA is radically different from the standard discretization based view of ADA, and it opens avenues for utilizing the vast efficient gradient-based algorithms available for continuous optimization problems. Results over multiple benchmarks demonstrate the efficiency improvement of this work compared with previous methods.

[Neural Routing by Memory](#)

- Kaipeng Zhang, Zhenqiang Li, Zhifeng Li, Wei Liu, Yoichi Sato
- abstract@[open-review](#): Recent Convolutional Neural Networks (CNNs) have achieved significant success by stacking multiple convolutional blocks, named procedures in this paper, to extract semantic features. However, they use the same procedure sequence for all inputs, regardless of the intermediate features. This paper proffers a simple yet effective idea of constructing parallel procedures and assigning similar intermediate features to the same specialized procedures in a divide-and-conquer fashion. It relieves each procedure's learning difficulty and thus leads to superior performance. Specifically, we propose a routing-by-memory mechanism for existing CNN architectures. In each stage of the network, we introduce parallel Procedural Units (PUs). A PU consists of a memory head and a procedure. The memory head maintains a summary of a type of features. For an intermediate feature, we search its closest memory and forward it to the corresponding procedure in both training and testing. In this way, different procedures are tailored to different features and therefore tackle them better. Networks with the proposed mechanism can be trained efficiently using a four-step training strategy. Experimental results show that our method improves VGGNet, ResNet, and EfficientNet's accuracies on Tiny ImageNet, ImageNet, and CIFAR-100 benchmarks with a negligible extra computational cost.

[GeoMol: Torsional Geometric Generation of Molecular 3D Conformer Ensembles](#)

- Octavian Ganea, Lagnajit Pattanaik, Connor Coley, Regina Barzilay, Klavs Jensen, William Green, Tommi Jaakkola
- abstract@[open-review](#): Prediction of a molecule's 3D conformer ensemble from the molecular graph holds a key role in areas of cheminformatics and drug discovery. Existing generative models have several drawbacks including lack of modeling important molecular geometry elements (e.g., torsion angles), separate optimization stages prone to error accumulation, and the need for structure fine-tuning based on approximate classical force-fields or computationally expensive methods. We propose GEOMOL --- an end-to-end, non-autoregressive, and SE(3)-invariant machine learning approach to generate distributions of low-energy molecular 3D conformers. Leveraging the power of message passing neural networks (MPNNs) to capture local and global graph information, we predict local atomic 3D structures and torsion angles, avoiding unnecessary over-parameterization of the geometric degrees

of freedom (e.g., one angle per non-terminal bond). Such local predictions suffice both for both the training loss computation and for the full deterministic conformer assembly (at test time). We devise a non-adversarial optimal transport based loss function to promote diverse conformer generation. GEOMOL predominantly outperforms popular open-source, commercial, or state-of-the-art machine learning (ML) models, while achieving significant speed-ups. We expect such differentiable 3D structure generators to significantly impact molecular modeling and related applications.

[CANITA: Faster Rates for Distributed Convex Optimization with Communication Compression](#)

- Zhize Li, Peter Richtarik
- abstract@[open-review](#): Due to the high communication cost in distributed and federated learning, methods relying on compressed communication are becoming increasingly popular. Besides, the best theoretically and practically performing gradient-type methods invariably rely on some form of acceleration/momentum to reduce the number of communications (faster convergence), e.g., Nesterov's accelerated gradient descent [31, 32] and Adam [14]. In order to combine the benefits of communication compression and convergence acceleration, we propose a \emph{compressed and accelerated} gradient method based on ANITA [20] for distributed optimization, which we call CANITA. Our CANITA achieves the \emph{first accelerated rate} $\mathcal{O}(\sqrt{\Big(1+\sqrt{\frac{\omega^3}{n}}\Big)\frac{L}{\epsilon}} + \omega\big(\frac{1}{\epsilon}\big)^{\frac{1}{3}})$, which improves upon the state-of-the-art non-accelerated rate $\mathcal{O}(\left(1+\frac{\omega}{n}\right)\frac{L}{\epsilon} + \frac{\omega^2+\omega}{\omega+n})$ of DIANA [12] for distributed general convex problems, where ϵ is the target error, L is the smooth parameter of the objective, n is the number of machines/devices, and ω is the compression parameter (larger ω means more compression can be applied, and no compression implies $\omega=0$). Our results show that as long as the number of devices n is large (often true in distributed/federated learning), or the compression ω is not very high, CANITA achieves the faster convergence rate $\mathcal{O}(\sqrt{\frac{L}{\epsilon}})$, i.e., the number of communication rounds is $\mathcal{O}(\sqrt{\frac{L}{\epsilon}})$ (vs. $\mathcal{O}(\frac{L}{\epsilon})$ achieved by previous works). As a result, CANITA enjoys the advantages of both compression (compressed communication in each round) and acceleration (much fewer communication rounds).

[Drop-DTW: Aligning Common Signal Between Sequences While Dropping Outliers](#)

- Mikita Dvornik, Isma Hadji, Konstantinos G. Derpanis, Animesh Garg, Allan Jepson
- abstract@[open-review](#): In this work, we consider the problem of sequence-to-sequence alignment for signals containing outliers. Assuming the absence of outliers, the standard Dynamic Time Warping (DTW) algorithm efficiently computes the optimal alignment between two (generally) variable-length sequences. While DTW is robust to temporal shifts and dilations of the signal, it fails to align sequences in a meaningful way in the presence of outliers that can be arbitrarily interspersed in the sequences. To address this problem, we introduce Drop-DTW, a novel algorithm that aligns the common signal between the sequences while automatically dropping the outlier elements from the matching. The entire procedure is implemented as a single dynamic program that is efficient and fully differentiable. In our experiments, we show that Drop-DTW is a robust similarity measure for sequence retrieval and demonstrate its effectiveness as a training loss on diverse applications. With Drop-DTW, we address temporal step localization on instructional videos, representation learning from noisy videos, and cross-modal representation learning for audio-visual retrieval and localization. In all applications, we take a weakly- or unsupervised approach and demonstrate state-of-the-art results under these settings.

[Safe Reinforcement Learning with Natural Language Constraints](#)

- Tsung-Yen Yang, Michael Y Hu, Yinlam Chow, Peter J Ramadge, Karthik Narasimhan
- abstract@[open-review](#): While safe reinforcement learning (RL) holds great promise for many practical applications like robotics or autonomous cars, current approaches require specifying constraints in mathematical form. Such specifications demand domain expertise, limiting the adoption of safe RL. In this paper, we propose learning to interpret natural language constraints for safe RL. To this end, we first introduce HAZARDWORLD, a new multi-task benchmark that requires an agent to optimize reward while not violating constraints specified in free-form text. We then develop an agent with a modular architecture that can interpret and adhere to such textual constraints while learning new tasks. Our model consists of (1) a constraint interpreter that encodes textual constraints into spatial and temporal representations of forbidden states, and (2) a policy network that uses these representations to produce a policy achieving minimal constraint violations during training. Across different domains in HAZARDWORLD, we show that our method achieves higher rewards (up to 11x) and fewer constraint violations (by 1.8x) compared to existing approaches. However, in terms of absolute performance, HAZARDWORLD still poses significant challenges for agents to learn efficiently, motivating the need for future work.

[Compositional Modeling of Nonlinear Dynamical Systems with ODE-based Random Features](#)

- Thomas McDonald, Mauricio Álvarez
- abstract@[open-review](#): Effectively modeling phenomena present in highly nonlinear dynamical systems whilst also accurately quantifying uncertainty is a challenging task, which often requires problem-specific techniques. We present a novel, domain-agnostic approach to tackling this problem, using compositions of physics-informed random features, derived from ordinary differential equations. The architecture of our model leverages recent advances in approximate inference for deep Gaussian processes, such as layer-wise weight-space approximations which allow us to incorporate random Fourier features, and stochastic variational inference for approximate Bayesian inference. We provide evidence that our model is capable of capturing highly nonlinear behaviour in real-world multivariate time series data. In addition, we find that our approach achieves comparable performance to a number of other probabilistic models on benchmark regression tasks.

[Implicit Semantic Response Alignment for Partial Domain Adaptation](#)

- Wenxiao Xiao, Zhengming Ding, Hongfu Liu
- abstract@[open-review](#): Partial Domain Adaptation (PDA) addresses the unsupervised domain adaptation problem where the target label space is a subset of the source label space. Most state-of-art PDA methods tackle the inconsistent label space by assigning weights to classes or individual samples, in an attempt to discard the source data that belongs to the irrelevant classes. However, we believe samples from those extra categories would still contain valuable information to promote positive transfer. In this paper, we propose the Implicit Semantic Response Alignment to explore the intrinsic relationships among different categories by applying a weighted schema on the feature level. Specifically, we design a class2vec module to extract the implicit semantic topics from the visual features. With an attention layer, we calculate the semantic response according to each implicit semantic topic. Then semantic responses of source and target data are aligned to retain the relevant information contained in multiple categories by weighting the features, instead of samples. Experiments on several cross-domain benchmark datasets demonstrate the effectiveness of our method over the state-of-the-art PDA methods. Moreover, we elaborate in-depth analyses to further explore implicit semantic alignment.

[ToAlign: Task-Oriented Alignment for Unsupervised Domain Adaptation](#)

- Guoqiang Wei, Cuiling Lan, Wenjun Zeng, Zhizheng Zhang, Zhibo Chen
- abstract@[open-review](#): Unsupervised domain adaptive classification intends to improve the classification performance on unlabeled target domain. To alleviate the adverse effect of domain shift, many approaches align the source and target domains in the feature space. However, a feature is usually taken as a whole for alignment without explicitly making domain alignment proactively serve the classification task, leading to sub-optimal solution. In this paper, we propose an effective Task-oriented Alignment (ToAlign) for unsupervised domain adaptation (UDA). We study what features should be aligned across domains and propose to make the domain alignment proactively serve classification by performing feature decomposition and alignment under the guidance of the prior knowledge induced from the classification task itself. Particularly, we explicitly decompose a feature in the source domain into a task-

related/discriminative feature that should be aligned, and a task-irrelevant feature that should be avoided/ignored, based on the classification meta-knowledge. Extensive experimental results on various benchmarks (e.g., Office-Home, Visda-2017, and DomainNet) under different domain adaptation settings demonstrate the effectiveness of ToAlign which helps achieve the state-of-the-art performance. The code is publicly available at <https://github.com/microsoft/UDA>.

[Prior-independent Dynamic Auctions for a Value-maximizing Buyer](#)

- Yuan Deng, Hanrui Zhang
- abstract@[open-review](#): We study prior-independent dynamic auction design with production costs for a value-maximizing buyer, a paradigm that is becoming prevalent recently following the development of automatic bidding algorithms in advertising platforms. In contrast to a utility-maximizing buyer, who maximizes the difference between her total value and total payment, a value-maximizing buyer aims to maximize her total value subject to a return on investment (ROI) constraint. Our main result is a dynamic mechanism with regret $\tilde{O}(T^{2/3})$, where T is the time horizon, against the first-best benchmark, i.e., the maximum amount of revenue the seller can extract assuming all values of the buyer are publicly known.

[Safe Reinforcement Learning by Imagining the Near Future](#)

- Garrett Thomas, Yuping Luo, Tengyu Ma
- abstract@[open-review](#): Safe reinforcement learning is a promising path toward applying reinforcement learning algorithms to real-world problems, where suboptimal behaviors may lead to actual negative consequences. In this work, we focus on the setting where unsafe states can be avoided by planning ahead a short time into the future. In this setting, a model-based agent with a sufficiently accurate model can avoid unsafe states. We devise a model-based algorithm that heavily penalizes unsafe trajectories, and derive guarantees that our algorithm can avoid unsafe states under certain assumptions. Experiments demonstrate that our algorithm can achieve competitive rewards with fewer safety violations in several continuous control tasks.

[Contrastive Active Inference](#)

- Pietro Mazzaglia, Tim Verbelen, Bart Dhoedt
- abstract@[open-review](#): Active inference is a unifying theory for perception and action resting upon the idea that the brain maintains an internal model of the world by minimizing free energy. From a behavioral perspective, active inference agents can be seen as self-evidencing beings that act to fulfill their optimistic predictions, namely preferred outcomes or goals. In contrast, reinforcement learning requires human-designed rewards to accomplish any desired outcome. Although active inference could provide a more natural self-supervised objective for control, its applicability has been limited because of the shortcomings in scaling the approach to complex environments. In this work, we propose a contrastive objective for active inference that strongly reduces the computational burden in learning the agent's generative model and planning future actions. Our method performs notably better than likelihood-based active inference in image-based tasks, while also being computationally cheaper and easier to train. We compare to reinforcement learning agents that have access to human-designed reward functions, showing that our approach closely matches their performance. Finally, we also show that contrastive methods perform significantly better in the case of distractors in the environment and that our method is able to generalize goals to variations in the background.

[Overparameterization Improves Robustness to Covariate Shift in High Dimensions](#)

- Nilesh Tripuraneni, Ben Adlam, Jeffrey Pennington
- abstract@[open-review](#): A significant obstacle in the development of robust machine learning models is covariate shift, a form of distribution shift that occurs when the input distributions of the training and test sets differ while the conditional label distributions remain the same. Despite the prevalence of covariate shift in real-world applications, a theoretical understanding in the context of modern machine learning has remained lacking. In this work, we examine the exact high-dimensional asymptotics of random feature regression under covariate shift and present a precise characterization of the limiting test error, bias, and variance in this setting. Our results motivate a natural partial order over covariate shifts that provides a sufficient condition for determining when the shift will harm (or even help) test performance. We find that overparameterized models exhibit enhanced robustness to covariate shift, providing one of the first theoretical explanations for this ubiquitous empirical phenomenon. Additionally, our analysis reveals an exact linear relationship between the in-distribution and out-of-distribution generalization performance, offering an explanation for this surprising recent observation.

[Logarithmic Regret in Feature-based Dynamic Pricing](#)

- Jianyu Xu, Yu-Xiang Wang
- abstract@[open-review](#): Feature-based dynamic pricing is an increasingly popular model of setting prices for highly differentiated products with applications in digital marketing, online sales, real estate and so on. The problem was formally studied as an online learning problem [Javanmard & Nazerzadeh, 2019] where a seller needs to propose prices on the fly for a sequence of T products based on their features x while having a small regret relative to the best --- "omniscient" --- pricing strategy she could have come up with in hindsight. We revisit this problem and provide two algorithms (EMLP and ONSP) for stochastic and adversarial feature settings, respectively, and prove the optimal $O(d\log{T})$ regret bounds for both. In comparison, the best existing results are $O(\min(\left(\frac{1}{\lambda}\right)^2 \log{T}, \sqrt{T}))$ and $O(T^{2/3})$ respectively, with λ being the smallest eigenvalue of $E[xx^T]$ that could be arbitrarily close to 0. We also prove an $\Omega(\sqrt{T})$ information-theoretic lower bound for a slightly more general setting, which demonstrates that "knowing-the-demand-curve" leads to an exponential improvement in feature-based dynamic pricing.

[Dimension-free empirical entropy estimation](#)

- Doron Cohen, Aryeh Kontorovich, Aaron Kolyk, Geoffrey Wolfer
- abstract@[open-review](#): We seek an entropy estimator for discrete distributions with fully empirical accuracy bounds. As stated, this goal is infeasible without some prior assumptions on the distribution. We discover that a certain information moment assumption renders the problem feasible. We argue that the moment assumption is natural and, in some sense, minimalistic --- weaker than finite support or tail decay conditions. Under the moment assumption, we provide the first finite-sample entropy estimates for infinite alphabets, nearly recovering the known minimax rates. Moreover, we demonstrate that our empirical bounds are significantly sharper than the state-of-the-art bounds, for various natural distributions and non-trivial sample regimes. Along the way, we give a dimension-free analogue of the Cover-Thomas result on entropy continuity (with respect to total variation distance) for finite alphabets, which may be of independent interest.

[Towards Biologically Plausible Convolutional Networks](#)

- Roman Pogodin, Yash Mehta, Timothy Lillicrap, Peter E Latham
- abstract@[open-review](#): Convolutional networks are ubiquitous in deep learning. They are particularly useful for images, as they reduce the number of parameters, reduce training time, and increase accuracy. However, as a model of the brain they are seriously problematic, since they require weight sharing - something real neurons simply cannot do. Consequently, while neurons in the brain can be locally connected (one of the features of convolutional networks), they cannot be convolutional. Locally connected but non-convolutional networks, however, significantly underperform convolutional ones. This is troublesome for studies that use convolutional networks to explain activity in the visual system. Here we study plausible

alternatives to weight sharing that aim at the same regularization principle, which is to make each neuron within a pool react similarly to identical inputs. The most natural way to do that is by showing the network multiple translations of the same image, akin to saccades in animal vision. However, this approach requires many translations, and doesn't remove the performance gap. We propose instead to add lateral connectivity to a locally connected network, and allow learning via Hebbian plasticity. This requires the network to pause occasionally for a sleep-like phase of "weight sharing". This method enables locally connected networks to achieve nearly convolutional performance on ImageNet and improves their fit to the ventral stream data, thus supporting convolutional networks as a model of the visual stream.

[DynamicViT: Efficient Vision Transformers with Dynamic Token Sparsification](#)

- Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, Cho-Jui Hsieh
- abstract@[open-review](#): Attention is sparse in vision transformers. We observe the final prediction in vision transformers is only based on a subset of most informative tokens, which is sufficient for accurate image recognition. Based on this observation, we propose a dynamic token sparsification framework to prune redundant tokens progressively and dynamically based on the input. Specifically, we devise a lightweight prediction module to estimate the importance score of each token given the current features. The module is added to different layers to prune redundant tokens hierarchically. To optimize the prediction module in an end-to-end manner, we propose an attention masking strategy to differentiably prune a token by blocking its interactions with other tokens. Benefiting from the nature of self-attention, the unstructured sparse tokens are still hardware friendly, which makes our framework easy to achieve actual speed-up. By hierarchically pruning 66% of the input tokens, our method greatly reduces 31% \$\sim\$ 37% FLOPs and improves the throughput by over 40% while the drop of accuracy is within 0.5% for various vision transformers. Equipped with the dynamic token sparsification framework, DynamicViT models can achieve very competitive complexity/accuracy trade-offs compared to state-of-the-art CNNs and vision transformers on ImageNet. Code is available at <https://github.com/raoyongming/DynamicViT>

[Learning Transferable Adversarial Perturbations](#)

- Krishna kanth Nakka, Mathieu Salzmann
- abstract@[open-review](#): While effective, deep neural networks (DNNs) are vulnerable to adversarial attacks. In particular, recent work has shown that such attacks could be generated by another deep network, leading to significant speedups over optimization-based perturbations. However, the ability of such generative methods to generalize to different test-time situations has not been systematically studied. In this paper, we, therefore, investigate the transferability of generated perturbations when the conditions at inference time differ from the training ones in terms of the target architecture, target data, and target task. Specifically, we identify the mid-level features extracted by the intermediate layers of DNNs as common ground across different architectures, datasets, and tasks. This lets us introduce a loss function based on such mid-level features to learn an effective, transferable perturbation generator. Our experiments demonstrate that our approach outperforms the state-of-the-art universal and transferable attack strategies.

[PortaSpeech: Portable and High-Quality Generative Text-to-Speech](#)

- Yi Ren, Jinglin Liu, Zhou Zhao
- abstract@[open-review](#): Non-autoregressive text-to-speech (NAR-TTS) models such as FastSpeech 2 and Glow-TTS can synthesize high-quality speech from the given text in parallel. After analyzing two kinds of generative NAR-TTS models (VAE and normalizing flow), we find that: VAE is good at capturing the long-range semantics features (e.g., prosody) even with small model size but suffers from blurry and unnatural results; and normalizing flow is good at reconstructing the frequency bin-wise details but performs poorly when the number of model parameters is limited. Inspired by these observations, to generate diverse speech with natural details and rich prosody using a lightweight architecture, we propose PortaSpeech, a portable and high-quality generative text-to-speech model. Specifically, 1) to model both the prosody and mel-spectrogram details accurately, we adopt a lightweight VAE with an enhanced prior followed by a flow-based post-net with strong conditional inputs as the main architecture. 2) To further compress the model size and memory footprint, we introduce the grouped parameter sharing mechanism to the affine coupling layers in the post-net. 3) To improve the expressiveness of synthesized speech and reduce the dependency on accurate fine-grained alignment between text and speech, we propose a linguistic encoder with mixture alignment combining hard word-level alignment and soft phoneme-level alignment, which explicitly extracts word-level semantic information. Experimental results show that PortaSpeech outperforms other TTS models in both voice quality and prosody modeling in terms of subjective and objective evaluation metrics, and shows only a slight performance degradation when reducing the model parameters to 6.7M (about 4x model size and 3x runtime memory compression ratio compared with FastSpeech 2). Our extensive ablation studies demonstrate that each design in PortaSpeech is effective.

[Exponential Graph is Provably Efficient for Decentralized Deep Training](#)

- Bicheng Ying, Kun Yuan, Yiming Chen, Hanbin Hu, PAN PAN, Wotao Yin
- abstract@[open-review](#): Decentralized SGD is an emerging training method for deep learning known for its much less (thus faster) communication per iteration, which relaxes the averaging step in parallel SGD to inexact averaging. The less exact the averaging is, however, the more the total iterations the training needs to take. Therefore, the key to making decentralized SGD efficient is to realize nearly-exact averaging using little communication. This requires a skillful choice of communication topology, which is an under-studied topic in decentralized optimization. In this paper, we study so-called exponential graphs where every node is connected to $O(\log(n))$ neighbors and n is the total number of nodes. This work proves such graphs can lead to both fast communication and effective averaging simultaneously. We also discover that a sequence of $\log(n)$ one-peer exponential graphs, in which each node communicates to one single neighbor per iteration, can together achieve exact averaging. This favorable property enables one-peer exponential graph to average as effective as its static counterpart but communicates more efficiently. We apply these exponential graphs in decentralized (momentum) SGD to obtain the state-of-the-art balance between per-iteration communication and iteration complexity among all commonly-used topologies. Experimental results on a variety of tasks and models demonstrate that decentralized (momentum) SGD over exponential graphs promises both fast and high-quality training. Our code is implemented through BlueFog and available at <https://github.com/Bluefog-Lib/NeurIPS2021-Exponential-Graph>.

[CLIP-It! Language-Guided Video Summarization](#)

- Medhini Narasimhan, Anna Rohrbach, Trevor Darrell
- abstract@[open-review](#): A generic video summary is an abridged version of a video that conveys the whole story and features the most important scenes. Yet the importance of scenes in a video is often subjective, and users should have the option of customizing the summary by using natural language to specify what is important to them. Further, existing models for fully automatic generic summarization have not exploited available language models, which can serve as an effective prior for saliency. This work introduces CLIP-It, a single framework for addressing both generic and query-focused video summarization, typically approached separately in the literature. We propose a language-guided multimodal transformer that learns to score frames in a video based on their importance relative to one another and their correlation with a user-defined query (for query-focused summarization) or an automatically generated dense video caption (for generic video summarization). Our model can be extended to the unsupervised setting by training without ground-truth supervision. We outperform baselines and prior work by a significant margin on both standard video summarization datasets (TVSum and SumMe) and a query-focused video summarization dataset (QFVS). Particularly, we achieve large improvements in the transfer setting, attesting to our method's strong generalization capabilities.

[Learning Treatment Effects in Panels with General Intervention Patterns](#)

- Vivek Farias, Andrew Li, Tianyi Peng

- abstract@[open-review](#): The problem of causal inference with panel data is a central econometric question. The following is a fundamental version of this problem: Let M^A be a low rank matrix and E be a zero-mean noise matrix. For a 'treatment' matrix Z with entries in $\{0,1\}$ we observe the matrix O with entries $O_{ij} := M^A_{ij} + E_{ij} + \mathcal{T}_{ij} Z_{ij}$ where \mathcal{T}_{ij} are unknown, heterogeneous treatment effects. The problem requires we estimate the average treatment effect $\tau := \sum_{ij} \mathcal{T}_{ij} Z_{ij} / \sum_{ij} Z_{ij}$. The synthetic control paradigm provides an approach to estimating τ when Z places support on a single row. This paper extends that framework to allow rate-optimal recovery of τ^* for general Z , thus broadly expanding its applicability. Our guarantees are the first of their type in this general setting. Computational experiments on synthetic and real-world data show a substantial advantage over competing estimators.

[Lossy Compression for Lossless Prediction](#)

- Yann Dubois, Benjamin Bloem-Reddy, Karen Ullrich, Chris J. Maddison
- abstract@[open-review](#): Most data is automatically collected and only ever "seen" by algorithms. Yet, data compressors preserve perceptual fidelity rather than just the information needed by algorithms performing downstream tasks. In this paper, we characterize the bit-rate required to ensure high performance on all predictive tasks that are invariant under a set of transformations, such as data augmentations. Based on our theory, we design unsupervised objectives for training neural compressors. Using these objectives, we train a generic image compressor that achieves substantial rate savings (more than 1000x on ImageNet) compared to JPEG on 8 datasets, without decreasing downstream classification performance.

[From Optimality to Robustness: Adaptive Re-Sampling Strategies in Stochastic Bandits](#)

- Dorian Baudry, Patrick Saux, Odalric-Ambrym Maillard
- abstract@[open-review](#): The stochastic multi-arm bandit problem has been extensively studied under standard assumptions on the arm's distribution (e.g bounded with known support, exponential family, etc). These assumptions are suitable for many real-world problems but sometimes they require knowledge (on tails for instance) that may not be precisely accessible to the practitioner, raising the question of the robustness of bandit algorithms to model misspecification. In this paper we study a generic `Dirichlet Sampling` (DS) algorithm, based on pairwise comparisons of empirical indices computed with `re-sampling` of the arms' observations and a data-dependent `exploration bonus`. We show that different variants of this strategy achieve provably optimal regret guarantees when the distributions are bounded and logarithmic regret for semi-bounded distributions with a mild quantile condition. We also show that a simple tuning achieve robustness with respect to a large class of unbounded distributions, at the cost of slightly worse than logarithmic asymptotic regret. We finally provide numerical experiments showing the merits of DS in a decision-making problem on synthetic agriculture data.

[CCVS: Context-aware Controllable Video Synthesis](#)

- Guillaume Le Moing, Jean Ponce, Cordelia Schmid
- abstract@[open-review](#): This presentation introduces a self-supervised learning approach to the synthesis of new video clips from old ones, with several new key elements for improved spatial resolution and realism: It conditions the synthesis process on contextual information for temporal continuity and ancillary information for fine control. The prediction model is doubly autoregressive, in the latent space of an autoencoder for forecasting, and in image space for updating contextual information, which is also used to enforce spatio-temporal consistency through a learnable optical flow module. Adversarial training of the autoencoder in the appearance and temporal domains is used to further improve the realism of its output. A quantizer inserted between the encoder and the transformer in charge of forecasting future frames in latent space (and its inverse inserted between the transformer and the decoder) adds even more flexibility by affording simple mechanisms for handling multimodal ancillary information for controlling the synthesis process (e.g., a few sample frames, an audio track, a trajectory in image space) and taking into account the intrinsically uncertain nature of the future by allowing multiple predictions. Experiments with an implementation of the proposed approach give very good qualitative and quantitative results on multiple tasks and standard benchmarks.

[An Online Riemannian PCA for Stochastic Canonical Correlation Analysis](#)

- Zihang Meng, Rudrasis Chakraborty, Vikas Singh
- abstract@[open-review](#): We present an efficient stochastic algorithm (RSG+) for canonical correlation analysis (CCA) using a reparametrization of the projection matrices. We show how this reparametrization (into structured matrices), simple in hindsight, directly presents an opportunity to repurpose/adjust mature techniques for numerical optimization on Riemannian manifolds. Our developments nicely complement existing methods for this problem which either require $O(d^3)$ time complexity per iteration with $O(\frac{1}{\sqrt{t}})$ convergence rate (where d is the dimensionality) or only extract the top 1 component with $O(\frac{1}{t})$ convergence rate. In contrast, our algorithm offers a strict improvement for this classical problem: it achieves $O(d^{2k})$ runtime complexity per iteration for extracting the top k canonical components with $O(\frac{1}{t})$ convergence rate. While the paper primarily focuses on the formulation and technical analysis of its properties, our experiments show that the empirical behavior on common datasets is quite promising. We also explore a potential application in training fair models where the label of protected attribute is missing or otherwise unavailable.

[Predify: Augmenting deep neural networks with brain-inspired predictive coding dynamics](#)

- Bhavin Choksi, Milad Mozafari, Callum Biggs O'May, B. ADOR, Andrea Alamia, Rufin VanRullen
- abstract@[open-review](#): Deep neural networks excel at image classification, but their performance is far less robust to input perturbations than human perception. In this work we explore whether this shortcoming may be partly addressed by incorporating brain-inspired recurrent dynamics in deep convolutional networks. We take inspiration from a popular framework in neuroscience: "predictive coding". At each layer of the hierarchical model, generative feedback "predicts" (i.e., reconstructs) the pattern of activity in the previous layer. The reconstruction errors are used to iteratively update the network's representations across timesteps, and to optimize the network's feedback weights over the natural image dataset--a form of unsupervised training. We show that implementing this strategy into two popular networks, VGG16 and EfficientNetB0, improves their robustness against various corruptions and adversarial attacks. We hypothesize that other feedforward networks could similarly benefit from the proposed framework. To promote research in this direction, we provide an open-sourced PyTorch-based package called `Predify`, which can be used to implement and investigate the impacts of the predictive coding dynamics in any convolutional neural network.

[Deep Extrapolation for Attribute-Enhanced Generation](#)

- Alvin Chan, Ali Madani, Ben Krause, Nikhil Naik
- abstract@[open-review](#): Attribute extrapolation in sample generation is challenging for deep neural networks operating beyond the training distribution. We formulate a new task for extrapolation in sequence generation, focusing on natural language and proteins, and propose GENhance, a generative framework that enhances attributes through a learned latent space. Trained on movie reviews and a computed protein stability dataset, GENhance can generate strongly-positive text reviews and highly stable protein sequences without being exposed to similar data during training. We release our benchmark tasks and models to contribute to the study of generative modeling extrapolation and data-driven design in biology and chemistry.

[Generalized Data Weighting via Class-Level Gradient Manipulation](#)

- Can Chen, Shuhao Zheng, Xi Chen, Erqun Dong, Xue (Steve) Liu, Hao Liu, Dejing Dou
- abstract@[open-review](#): Label noise and class imbalance are two major issues coexisting in real-world datasets. To alleviate the two issues, state-of-the-art methods reweight each instance by leveraging a small amount of clean and unbiased data. Yet, these methods overlook class-level information within each instance, which can be further utilized to improve performance. To this end, in this paper, we propose Generalized Data Weighting (GDW) to simultaneously mitigate label noise and class imbalance by manipulating gradients at the class level. To be specific, GDW unrolls the loss gradient to class-level gradients by the chain rule and reweights the flow of each gradient separately. In this way, GDW achieves remarkable performance improvement on both issues. Aside from the performance gain, GDW efficiently obtains class-level weights without introducing any extra computational cost compared with instance weighting methods. Specifically, GDW performs a gradient descent step on class-level weights, which only relies on intermediate gradients. Extensive experiments in various settings verify the effectiveness of GDW. For example, GDW outperforms state-of-the-art methods by \$2.56\%\$ under the \$60\%\$ uniform noise setting in CIFAR10. Our code is available at <https://github.com/GGchen1997/GDW-NIPS2021>.

[Slow Learning and Fast Inference: Efficient Graph Similarity Computation via Knowledge Distillation](#)

- Can Qin, Handong Zhao, Lichen Wang, Huan Wang, Yulun Zhang, Yun Fu
- abstract@[open-review](#): Graph Similarity Computation (GSC) is essential to wide-ranging graph applications such as retrieval, plagiarism/anomaly detection, etc. The exact computation of graph similarity, e.g., Graph Edit Distance (GED), is an NP-hard problem that cannot be exactly solved within an adequate time given large graphs. Thanks to the strong representation power of graph neural network (GNN), a variety of GNN-based inexact methods emerged. To capture the subtle difference across graphs, the key success is designing the dense interaction with features fusion at the early stage, which, however, is a trade-off between speed and accuracy. For slow learning of graph similarity, this paper proposes a novel early-fusion approach by designing a co-attention-based feature fusion network on multilevel GNN features. To further improve the speed without much accuracy drop, we introduce an efficient GSC solution by distilling the knowledge from the slow early-fusion model to the student one for fast inference. Such a student model also enables the offline collection of individual graph embeddings, speeding up the inference time in orders. To address the instability through knowledge transfer, we decompose the dynamic joint embedding into the static pseudo individual ones for precise teacher-student alignment. The experimental analysis on the real-world datasets demonstrates the superiority of our approach over the state-of-the-art methods on both accuracy and efficiency. Particularly, we speed up the prior art by more than 10x on the benchmark AIDS data.

[Meta Learning Backpropagation And Improving It](#)

- Louis Kirsch, Jürgen Schmidhuber
- abstract@[open-review](#): Many concepts have been proposed for meta learning with neural networks (NNs), e.g., NNs that learn to reprogram fast weights, Hebbian plasticity, learned learning rules, and meta recurrent NNs. Our Variable Shared Meta Learning (VSML) unifies the above and demonstrates that simple weight-sharing and sparsity in an NN is sufficient to express powerful learning algorithms (LAs) in a reusable fashion. A simple implementation of VSML where the weights of a neural network are replaced by tiny LSTMs allows for implementing the backpropagation LA solely by running in forward-mode. It can even meta learn new LAs that differ from online backpropagation and generalize to datasets outside of the meta training distribution without explicit gradient calculation. Introspection reveals that our meta learned LAs learn through fast association in a way that is qualitatively different from gradient descent.

[Posterior Meta-Replay for Continual Learning](#)

- Christian Henning, Maria Cervera, Francesco D'Angelo, Johannes von Oswald, Regina Traber, Benjamin Ehret, Seijin Kobayashi, Benjamin F. Grewe, João Sacramento
- abstract@[open-review](#): Learning a sequence of tasks without access to i.i.d. observations is a widely studied form of continual learning (CL) that remains challenging. In principle, Bayesian learning directly applies to this setting, since recursive and one-off Bayesian updates yield the same result. In practice, however, recursive updating often leads to poor trade-off solutions across tasks because approximate inference is necessary for most models of interest. Here, we describe an alternative Bayesian approach where task-conditioned parameter distributions are continually inferred from data. We offer a practical deep learning implementation of our framework based on probabilistic task-conditioned hypernetworks, an approach we term posterior meta-replay. Experiments on standard benchmarks show that our probabilistic hypernetworks compress sequences of posterior parameter distributions with virtually no forgetting. We obtain considerable performance gains compared to existing Bayesian CL methods, and identify task inference as our major limiting factor. This limitation has several causes that are independent of the considered sequential setting, opening up new avenues for progress in CL.

[Optimizing Reusable Knowledge for Continual Learning via Metalearning](#)

- Julio Hurtado, Alain Raymond, Alvaro Soto
- abstract@[open-review](#): When learning tasks over time, artificial neural networks suffer from a problem known as Catastrophic Forgetting (CF). This happens when the weights of a network are overwritten during the training of a new task causing forgetting of old information. To address this issue, we propose MetA Reusable Knowledge or MARK, a new method that fosters weight reusability instead of overwriting when learning a new task. Specifically, MARK keeps a set of shared weights among tasks. We envision these shared weights as a common Knowledge Base (KB) that is not only used to learn new tasks, but also enriched with new knowledge as the model learns new tasks. Key components behind MARK are two-fold. On the one hand, a metalearning approach provides the key mechanism to incrementally enrich the KB with new knowledge and to foster weight reusability among tasks. On the other hand, a set of trainable masks provides the key mechanism to selectively choose from the KB relevant weights to solve each task. By using MARK, we achieve state of the art results in several popular benchmarks, surpassing the best performing methods in terms of average accuracy by over 10% on the 20-Split-MiniImageNet dataset, while achieving almost zero forgetfulness using 55% of the number of parameters. Furthermore, an ablation study provides evidence that, indeed, MARK is learning reusable knowledge that is selectively used by each task.

[A sampling-based circuit for optimal decision making](#)

- Camille Rullán Buxó, Cristina Savin
- abstract@[open-review](#): Many features of human and animal behavior can be understood in the framework of Bayesian inference and optimal decision making, but the biological substrate of such processes is not fully understood. Neural sampling provides a flexible code for probabilistic inference in high dimensions and explains key features of sensory responses under experimental manipulations of uncertainty. However, since it encodes uncertainty implicitly, across time and neurons, it remains unclear how such representations can be used for decision making. Here we propose a spiking network model that maps neural samples of a task-specific marginal distribution into an instantaneous representation of uncertainty via a procedure inspired by online kernel density estimation, so that its output can be readily used for decision making. Our model is consistent with experimental results at the level of single neurons and populations, and makes predictions for how neural responses and decisions could be modulated by uncertainty and prior biases. More generally, our work brings together conflicting perspectives on probabilistic brain computation.

[Compressed Video Contrastive Learning](#)

- Yuqi Huo, Mingyu Ding, Haoyu Lu, Nanyi Fei, Zhiwu Lu, Ji-Rong Wen, Ping Luo
- abstract@[open-review](#): This work concerns self-supervised video representation learning (SSVRL), one topic that has received much attention recently. Since videos are storage-intensive and contain a rich source of visual content, models designed for SSVRL are expected to be storage- and computation-efficient, as well as effective. However, most existing methods only focus on one of the two objectives, failing to consider both at the same time. In this

work, for the first time, the seemingly contradictory goals are simultaneously achieved by exploiting compressed videos and capturing mutual information between two input streams. Specifically, a novel Motion Vector based Cross Guidance Contrastive learning approach (MVCGC) is proposed. For storage and computation efficiency, we choose to directly decode RGB frames and motion vectors (that resemble low-resolution optical flows) from compressed videos on-the-fly. To enhance the representation ability of the motion vectors, hence the effectiveness of our method, we design a cross guidance contrastive learning algorithm based on multi-instance InfoNCE loss, where motion vectors can take supervision signals from RGB frames and vice versa. Comprehensive experiments on two downstream tasks show that our MVCGC yields new state-of-the-art while being significantly more efficient than its competitors.

[Uniform-PAC Bounds for Reinforcement Learning with Linear Function Approximation](#)

- Jiafan He, Dongruo Zhou, Quanquan Gu
- abstract@[open-review](#): We study reinforcement learning (RL) with linear function approximation. Existing algorithms for this problem only have high-probability regret and/or Probably Approximately Correct (PAC) sample complexity guarantees, which cannot guarantee the convergence to the optimal policy. In this paper, in order to overcome the limitation of existing algorithms, we propose a new algorithm called FLUTE, which enjoys uniform-PAC convergence to the optimal policy with high probability. The uniform-PAC guarantee is the strongest possible guarantee for reinforcement learning in the literature, which can directly imply both PAC and high probability regret bounds, making our algorithm superior to all existing algorithms with linear function approximation. At the core of our algorithm is a novel minimax value function estimator and a multi-level partition scheme to select the training samples from historical observations. Both of these techniques are new and of independent interest.

[Attention Bottlenecks for Multimodal Fusion](#)

- Arsha Nagrani, Shan Yang, Anurag Arnab, Aren Jansen, Cordelia Schmid, Chen Sun
- abstract@[open-review](#): Humans perceive the world by concurrently processing and fusing high-dimensional inputs from multiple modalities such as vision and audio. Machine perception models, in stark contrast, are typically modality-specific and optimised for unimodal benchmarks. A common approach for building multimodal models is to simply combine multiple of these modality-specific architectures using late-stage fusion of final representations or predictions ('late-fusion'). Instead, we introduce a novel transformer based architecture that uses 'attention bottlenecks' for modality fusion at multiple layers. Compared to traditional pairwise self-attention, these bottlenecks force information between different modalities to pass through a small number of 'bottleneck' latent units, requiring the model to collate and condense the most relevant information in each modality and only share what is necessary. We find that such a strategy improves fusion performance, at the same time reducing computational cost. We conduct thorough ablation studies, and achieve state-of-the-art results on multiple audio-visual classification benchmarks including AudioSet, Epic-Kitchens and VGGSound. All code and models will be released.

[Convergence of adaptive algorithms for constrained weakly convex optimization](#)

- Ahmet Alacaoglu, Yura Malitsky, Volkan Cevher
- abstract@[open-review](#): We analyze the adaptive first order algorithm AMSGrad, for solving a constrained stochastic optimization problem with a weakly convex objective. We prove the $\tilde{O}(t^{-1/2})$ rate of convergence for the squared norm of the gradient of Moreau envelope, which is the standard stationarity measure for this class of problems. It matches the known rates that adaptive algorithms enjoy for the specific case of unconstrained smooth nonconvex stochastic optimization. Our analysis works with mini-batch size of 1 , constant first and second order moment parameters, and possibly unbounded optimization domains. Finally, we illustrate the applications and extensions of our results to specific problems and algorithms.

[On the Convergence of Step Decay Step-Size for Stochastic Optimization](#)

- Xiaoyu Wang, Sindri Magnússon, Mikael Johansson
- abstract@[open-review](#): The convergence of stochastic gradient descent is highly dependent on the step-size, especially on non-convex problems such as neural network training. Step decay step-size schedules (constant and then cut) are widely used in practice because of their excellent convergence and generalization qualities, but their theoretical properties are not yet well understood. We provide convergence results for step decay in the non-convex regime, ensuring that the gradient norm vanishes at an $O(\ln T/\sqrt{T})$ rate. We also provide near-optimal (and sometimes provably tight) convergence guarantees for general, possibly non-smooth, convex and strongly convex problems. The practical efficiency of the step decay step-size is demonstrated in several large-scale deep neural network training tasks.

[BernNet: Learning Arbitrary Graph Spectral Filters via Bernstein Approximation](#)

- Mingguo He, Zhewei Wei, zengfeng Huang, Hongteng Xu
- abstract@[open-review](#): Many representative graph neural networks, e.g., GPR-GNN and ChebNet, approximate graph convolutions with graph spectral filters. However, existing work either applies predefined filter weights or learns them without necessary constraints, which may lead to oversimplified or ill-posed filters. To overcome these issues, we propose BernNet, a novel graph neural network with theoretical support that provides a simple but effective scheme for designing and learning arbitrary graph spectral filters. In particular, for any filter over the normalized Laplacian spectrum of a graph, our BernNet estimates it by an order- K Bernstein polynomial approximation and designs its spectral property by setting the coefficients of the Bernstein basis. Moreover, we can learn the coefficients (and the corresponding filter weights) based on observed graphs and their associated signals and thus achieve the BernNet specialized for the data. Our experiments demonstrate that BernNet can learn arbitrary spectral filters, including complicated band-rejection and comb filters, and it achieves superior performance in real-world graph modeling tasks. Code is available at <https://github.com/ivam-he/BernNet>.

[Co-evolution Transformer for Protein Contact Prediction](#)

- He Zhang, Fusong Ju, Jianwei Zhu, Liang He, Bin Shao, Nanning Zheng, Tie-Yan Liu
- abstract@[open-review](#): Proteins are the main machinery of life and protein functions are largely determined by their 3D structures. The measurement of the pairwise proximity between amino acids of a protein, known as inter-residue contact map, well characterizes the structural information of a protein. Protein contact prediction (PCP) is an essential building block of many protein structure related applications. The prevalent approach to contact prediction is based on estimating the inter-residue contacts using hand-crafted coevolutionary features derived from multiple sequence alignments (MSAs). To mitigate the information loss caused by hand-crafted features, some recently proposed methods try to learn residue co-evolutions directly from MSAs. These methods generally derive coevolutionary features by aggregating the learned residue representations from individual sequences with equal weights, which is inconsistent with the premise that residue co-evolutions are a reflection of collective covariation patterns of numerous homologous proteins. Moreover, non-homologous residues and gaps commonly exist in MSAs. By aggregating features from all homologs equally, the non-homologous information may cause misestimation of the residue co-evolutions. To overcome these issues, we propose an attention-based architecture, Co-evolution Transformer (CoT), for PCP. CoT jointly considers the information from all homologous sequences in the MSA to better capture global coevolutionary patterns. To mitigate the influence of the non-homologous information, CoT selectively aggregates the features from different homologs by assigning smaller weights to non-homologous sequences or residue pairs. Extensive experiments on two rigorous benchmark datasets demonstrate the effectiveness of CoT. In particular, CoT achieves a 51.6% top-L long-range precision score for the Free Modeling (FM) domains on the CASP14 benchmark, which outperforms the winner group of CASP14 contact prediction challenge by 9.8% .

Unsupervised Foreground Extraction via Deep Region Competition

- Peiyu Yu, Sirui Xie, Xiaojian Ma, Yixin Zhu, Ying Nian Wu, Song-Chun Zhu
- abstract@[open-review](#): We present Deep Region Competition (DRC), an algorithm designed to extract foreground objects from images in a fully unsupervised manner. Foreground extraction can be viewed as a special case of generic image segmentation that focuses on identifying and disentangling objects from the background. In this work, we rethink the foreground extraction by reconciling energy-based prior with generative image modeling in the form of Mixture of Experts (MoE), where we further introduce the learned pixel re-assignment as the essential inductive bias to capture the regularities of background regions. With this modeling, the foreground-background partition can be naturally found through Expectation-Maximization (EM). We show that the proposed method effectively exploits the interaction between the mixture components during the partitioning process, which closely connects to region competition, a seminal approach for generic image segmentation. Experiments demonstrate that DRC exhibits more competitive performances on complex real-world data and challenging multi-object scenes compared with prior methods. Moreover, we show empirically that DRC can potentially generalize to novel foreground objects even from categories unseen during training.

Leveraging Spatial and Temporal Correlations in Sparsified Mean Estimation

- Divyansh Jhunjhunwala, Ankur Mallick, Advait Gadhikar, Swanand Kadhe, Gauri Joshi
- abstract@[open-review](#): We study the problem of estimating at a central server the mean of a set of vectors distributed across several nodes (one vector per node). When the vectors are high-dimensional, the communication cost of sending entire vectors may be prohibitive, and it may be imperative for them to use sparsification techniques. While most existing work on sparsified mean estimation is agnostic to the characteristics of the data vectors, in many practical applications such as federated learning, there may be spatial correlations (similarities in the vectors sent by different nodes) or temporal correlations (similarities in the data sent by a single node over different iterations of the algorithm) in the data vectors. We leverage these correlations by simply modifying the decoding method used by the server to estimate the mean. We provide an analysis of the resulting estimation error as well as experiments for PCA, K-Means and Logistic Regression, which show that our estimators consistently outperform more sophisticated and expensive sparsification methods.

Last-iterate Convergence in Extensive-Form Games

- Chung-Wei Lee, Christian Kroer, Haipeng Luo
- abstract@[open-review](#): Regret-based algorithms are highly efficient at finding approximate Nash equilibria in sequential games such as poker games. However, most regret-based algorithms, including counterfactual regret minimization (CFR) and its variants, rely on iterate averaging to achieve convergence. Inspired by recent advances on last-iterate convergence of optimistic algorithms in zero-sum normal-form games, we study this phenomenon in sequential games, and provide a comprehensive study of last-iterate convergence for zero-sum extensive-form games with perfect recall (EFGs), using various optimistic regret-minimization algorithms over treeplexes. This includes algorithms using the vanilla entropy or squared Euclidean norm regularizers, as well as their dilated versions which admit more efficient implementation. In contrast to CFR, we show that all of these algorithms enjoy last-iterate convergence, with some of them even converging exponentially fast. We also provide experiments to further support our theoretical results.

Class-Incremental Learning via Dual Augmentation

- Fei Zhu, Zhen Cheng, Xu-yao Zhang, Cheng-lin Liu
- abstract@[open-review](#): Deep learning systems typically suffer from catastrophic forgetting of past knowledge when acquiring new skills continually. In this paper, we emphasize two dilemmas, representation bias and classifier bias in class-incremental learning, and present a simple and novel approach that employs explicit class augmentation (classAug) and implicit semantic augmentation (semanAug) to address the two biases, respectively. On the one hand, we propose to address the representation bias by learning transferable and diverse representations. Specifically, we investigate the feature representations in incremental learning based on spectral analysis and present a simple technique called classAug, to let the model see more classes during training for learning representations transferable across classes. On the other hand, to overcome the classifier bias, semanAug implicitly involves the simultaneous generating of an infinite number of instances of old classes in the deep feature space, which poses tighter constraints to maintain the decision boundary of previously learned classes. Without storing any old samples, our method can perform comparably with representative data replay based approaches.

Robust and Fully-Dynamic Coreset for Continuous-and-Bounded Learning (With Outliers) Problems

- Zixiu Wang, Yiwen Guo, Hu Ding
- abstract@[open-review](#): In many machine learning tasks, a common approach for dealing with large-scale data is to build a small summary, {\em e.g.,} coreset, that can efficiently represent the original input. However, real-world datasets usually contain outliers and most existing coreset construction methods are not resilient against outliers (in particular, an outlier can be located arbitrarily in the space by an adversarial attacker). In this paper, we propose a novel robust coreset method for the {\em continuous-and-bounded learning} problems (with outliers) which includes a broad range of popular optimization objectives in machine learning, {\em e.g.,} logistic regression and \$k\$-means clustering. Moreover, our robust coreset can be efficiently maintained in fully-dynamic environment. To the best of our knowledge, this is the first robust and fully-dynamic coreset construction method for these optimization problems. Another highlight is that our coreset size can depend on the doubling dimension of the parameter space, rather than the VC dimension of the objective function which could be very large or even challenging to compute. Finally, we conduct the experiments on real-world datasets to evaluate the effectiveness of our proposed robust coreset method.

Rethinking and Reweighting the Univariate Losses for Multi-Label Ranking: Consistency and Generalization

- Guoqiang Wu, Chongxuan LI, Kun Xu, Jun Zhu
- abstract@[open-review](#): The (partial) ranking loss is a commonly used evaluation measure for multi-label classification, which is usually optimized with convex surrogates for computational efficiency. Prior theoretical efforts on multi-label ranking mainly focus on (Fisher) consistency analyses. However, there is a gap between existing theory and practice --- some inconsistent pairwise losses can lead to promising performance, while some consistent univariate losses usually have no clear superiority in practice. To take a step towards filling up this gap, this paper presents a systematic study from two complementary perspectives of consistency and generalization error bounds of learning algorithms. We theoretically find two key factors of the distribution (or dataset) that affect the learning guarantees of algorithms: the instance-wise class imbalance and the label size \$c\$. Specifically, in an extremely imbalanced case, the algorithm with the consistent univariate loss has an error bound of \$O(c)\$, while the one with the inconsistent pairwise loss depends on \$O(\sqrt{c})\$ as shown in prior work. This may shed light on the superior performance of pairwise methods in practice, where real datasets are usually highly imbalanced. Moreover, we present an inconsistent reweighted univariate loss-based algorithm that enjoys an error bound of \$O(\sqrt{c})\$ for promising performance as well as the computational efficiency of univariate losses. Finally, experimental results confirm our theoretical findings.

Fair Clustering Under a Bounded Cost

- Seyed Esmaeili, Brian Brubach, Aravind Srinivasan, John Dickerson
- abstract@[open-review](#): Clustering is a fundamental unsupervised learning problem where a dataset is partitioned into clusters that consist of nearby points in a metric space. A recent variant, fair clustering, associates a color with each point representing its group membership and requires that each color has

(approximately) equal representation in each cluster to satisfy group fairness. In this model, the cost of the clustering objective increases due to enforcing fairness in the algorithm. The relative increase in the cost, the "price of fairness," can indeed be unbounded. Therefore, in this paper we propose to treat an upper bound on the clustering objective as a constraint on the clustering problem, and to maximize equality of representation subject to it. We consider two fairness objectives: the group utilitarian objective and the group egalitarian objective, as well as the group leximin objective which generalizes the group egalitarian objective. We derive fundamental lower bounds on the approximation of the utilitarian and egalitarian objectives and introduce algorithms with provable guarantees for them. For the leximin objective we introduce an effective heuristic algorithm. We further derive impossibility results for other natural fairness objectives. We conclude with experimental results on real-world datasets that demonstrate the validity of our algorithms.

Improving Calibration through the Relationship with Adversarial Robustness

- Yao Qin, Xuezhi Wang, Alex Beutel, Ed Chi
- abstract@[open-review](#): Neural networks lack adversarial robustness, i.e., they are vulnerable to adversarial examples that through small perturbations to inputs cause incorrect predictions. Further, trust is undermined when models give miscalibrated predictions, i.e., the predicted probability is not a good indicator of how much we should trust our model. In this paper, we study the connection between adversarial robustness and calibration and find that the inputs for which the model is sensitive to small perturbations (are easily attacked) are more likely to have poorly calibrated predictions. Based on this insight, we examine if calibration can be improved by addressing those adversarially unrobust inputs. To this end, we propose Adversarial Robustness based Adaptive Label Smoothing (AR-AdaLS) that integrates the correlations of adversarial robustness and calibration into training by adaptively softening labels for an example based on how easily it can be attacked by an adversary. We find that our method, taking the adversarial robustness of the in-distribution data into consideration, leads to better calibration over the model even under distributional shifts. In addition, AR-AdaLS can also be applied to an ensemble model to further improve model calibration.

Credal Self-Supervised Learning

- Julian Lienen, Eyke Hüllermeier
- abstract@[open-review](#): Self-training is an effective approach to semi-supervised learning. The key idea is to let the learner itself iteratively generate "pseudo-supervision" for unlabeled instances based on its current hypothesis. In combination with consistency regularization, pseudo-labeling has shown promising performance in various domains, for example in computer vision. To account for the hypothetical nature of the pseudo-labels, these are commonly provided in the form of probability distributions. Still, one may argue that even a probability distribution represents an excessive level of informedness, as it suggests that the learner precisely knows the ground-truth conditional probabilities. In our approach, we therefore allow the learner to label instances in the form of credal sets, that is, sets of (candidate) probability distributions. Thanks to this increased expressiveness, the learner is able to represent uncertainty and a lack of knowledge in a more flexible and more faithful manner. To learn from weakly labeled data of that kind, we leverage methods that have recently been proposed in the realm of so-called superset learning. In an exhaustive empirical evaluation, we compare our methodology to state-of-the-art self-supervision approaches, showing competitive to superior performance especially in low-label scenarios incorporating a high degree of uncertainty.

Spot the Difference: Detection of Topological Changes via Geometric Alignment

- Per Steffen Czolbe, Aasa Feragen, Oswin Krause
- abstract@[open-review](#): Geometric alignment appears in a variety of applications, ranging from domain adaptation, optimal transport, and normalizing flows in machine learning; optical flow and learned augmentation in computer vision and deformable registration within biomedical imaging. A recurring challenge is the alignment of domains whose topology is not the same; a problem that is routinely ignored, potentially introducing bias in downstream analysis. As a first step towards solving such alignment problems, we propose an unsupervised algorithm for the detection of changes in image topology. The model is based on a conditional variational auto-encoder and detects topological changes between two images during the registration step. We account for both topological changes in the image under spatial variation and unexpected transformations. Our approach is validated on two tasks and datasets: detection of topological changes in microscopy images of cells, and unsupervised anomaly detection brain imaging.

Rethinking the Variational Interpretation of Accelerated Optimization Methods

- Peiyuan Zhang, Antonio Orvieto, Hadi Daneshmand
- abstract@[open-review](#): The continuous-time model of Nesterov's momentum provides a thought-provoking perspective for understanding the nature of the acceleration phenomenon in convex optimization. One of the main ideas in this line of research comes from the field of classical mechanics and proposes to link Nesterov's trajectory to the solution of a set of Euler-Lagrange equations relative to the so-called Bregman Lagrangian. In the last years, this approach led to the discovery of many new (stochastic) accelerated algorithms and provided a solid theoretical foundation for the design of structure-preserving accelerated methods. In this work, we revisit this idea and provide an in-depth analysis of the action relative to the Bregman Lagrangian from the point of view of calculus of variations. Our main finding is that, while Nesterov's method is a stationary point for the action, it is often not a minimizer but instead a saddle point for this functional in the space of differentiable curves. This finding challenges the main intuition behind the variational interpretation of Nesterov's method and provides additional insights into the intriguing geometry of accelerated paths.

Linear and Kernel Classification in the Streaming Model: Improved Bounds for Heavy Hitters

- Arvind Mahankali, David Woodruff
- abstract@[open-review](#): We study linear and kernel classification in the streaming model. For linear classification, we improve upon the algorithm of (Tai, et al. 2018), which solves the ℓ_1 point query problem on the optimal weight vector $w \in \mathbb{R}^d$ in sublinear space. We first give an algorithm solving the more difficult ℓ_2 point query problem on w , also in sublinear space. We also give an algorithm which solves the ℓ_2 heavy hitter problem on w , in sublinear space and running time. Finally, we give an algorithm which can deterministically solve the ℓ_1 point query problem on w , with sublinear space improving upon that of (Tai, et al. 2018). For kernel classification, if $w \in \mathbb{R}^d$ is the optimal weight vector classifying points in the stream according to their p -degree polynomial kernel, then we give an algorithm solving the ℓ_2 point query problem on w in $\text{poly}(\frac{p}{\log d}/\epsilon)$ space, and an algorithm solving the ℓ_2 heavy hitter problem in $\text{poly}(\frac{p}{\log d}/\epsilon)$ space and running time. Note that our space and running time are polynomial in p , making our algorithms well-suited to high-degree polynomial kernels and the Gaussian kernel (approximated by the polynomial kernel of degree $p = \Theta(\log T)$).

A PAC-Bayes Analysis of Adversarial Robustness

- Paul Viallard, Eric Guillaume VIDOT, Amaury Habrard, Emilie Morvant
- abstract@[open-review](#): We propose the first general PAC-Bayesian generalization bounds for adversarial robustness, that estimate, at test time, how much a model will be invariant to imperceptible perturbations in the input. Instead of deriving a worst-case analysis of the risk of a hypothesis over all the possible perturbations, we leverage the PAC-Bayesian framework to bound the averaged risk on the perturbations for majority votes (over the whole class of hypotheses). Our theoretically founded analysis has the advantage to provide general bounds (i) that are valid for any kind of attacks (i.e., the adversarial attacks), (ii) that are tight thanks to the PAC-Bayesian framework, (iii) that can be directly minimized during the learning phase to obtain a robust model on different attacks at test time.

SE(3)-equivariant prediction of molecular wavefunctions and electronic densities

- Oliver Unke, Mihail Bogojeski, Michael Gastegger, Mario Geiger, Tess Smidt, Klaus-Robert Müller
- abstract@[open-review](#): Machine learning has enabled the prediction of quantum chemical properties with high accuracy and efficiency, allowing to bypass computationally costly ab initio calculations. Instead of training on a fixed set of properties, more recent approaches attempt to learn the electronic wavefunction (or density) as a central quantity of atomistic systems, from which all other observables can be derived. This is complicated by the fact that wavefunctions transform non-trivially under molecular rotations, which makes them a challenging prediction target. To solve this issue, we introduce general SE(3)-equivariant operations and building blocks for constructing deep learning architectures for geometric point cloud data and apply them to reconstruct wavefunctions of atomistic systems with unprecedented accuracy. Our model achieves speedups of over three orders of magnitude compared to ab initio methods and reduces prediction errors by up to two orders of magnitude compared to the previous state-of-the-art. This accuracy makes it possible to derive properties such as energies and forces directly from the wavefunction in an end-to-end manner. We demonstrate the potential of our approach in a transfer learning application, where a model trained on low accuracy reference wavefunctions implicitly learns to correct for electronic many-body interactions from observables computed at a higher level of theory. Such machine-learned wavefunction surrogates pave the way towards novel semi-empirical methods, offering resolution at an electronic level while drastically decreasing computational cost. Additionally, the predicted wavefunctions can serve as initial guess in conventional ab initio methods, decreasing the number of iterations required to arrive at a converged solution, thus leading to significant speedups without any loss of accuracy or robustness. While we focus on physics applications in this contribution, the proposed equivariant framework for deep learning on point clouds is promising also beyond, say, in computer vision or graphics.

Modified Frank Wolfe in Probability Space

- Carson Kent, Jiajin Li, Jose Blanchet, Peter W Glynn
- abstract@[open-review](#): We propose a novel Frank-Wolfe (FW) procedure for the optimization of infinite-dimensional functionals of probability measures - a task which arises naturally in a wide range of areas including statistical learning (e.g. variational inference) and artificial intelligence (e.g. generative adversarial networks). Our FW procedure takes advantage of Wasserstein gradient flows and strong duality results recently developed in Distributionally Robust Optimization so that gradient steps (in the Wasserstein space) can be efficiently computed using finite-dimensional, convex optimization methods. We show how to choose the step sizes in order to guarantee exponentially fast iteration convergence, under mild assumptions on the functional to optimize. We apply our algorithm to a range of functionals arising from applications in nonparametric estimation.

Bayesian Optimization of Function Networks

- Raul Astudillo, Peter Frazier
- abstract@[open-review](#): We consider Bayesian optimization of the output of a network of functions, where each function takes as input the output of its parent nodes, and where the network takes significant time to evaluate. Such problems arise, for example, in reinforcement learning, engineering design, and manufacturing. While the standard Bayesian optimization approach observes only the final output, our approach delivers greater query efficiency by leveraging information that the former ignores: intermediate output within the network. This is achieved by modeling the nodes of the network using Gaussian processes and choosing the points to evaluate using, as our acquisition function, the expected improvement computed with respect to the implied posterior on the objective. Although the non-Gaussian nature of this posterior prevents computing our acquisition function in closed form, we show that it can be efficiently maximized via sample average approximation. In addition, we prove that our method is asymptotically consistent, meaning that it finds a globally optimal solution as the number of evaluations grows to infinity, thus generalizing previously known convergence results for the expected improvement. Notably, this holds even though our method might not evaluate the domain densely, instead leveraging problem structure to leave regions unexplored. Finally, we show that our approach dramatically outperforms standard Bayesian optimization methods in several synthetic and real-world problems.

Look at What I'm Doing: Self-Supervised Spatial Grounding of Narrations in Instructional Videos

- Reuben Tan, Bryan Plummer, Kate Saenko, Hailin Jin, Bryan Russell
- abstract@[open-review](#): We introduce the task of spatially localizing narrated interactions in videos. Key to our approach is the ability to learn to spatially localize interactions with self-supervision on a large corpus of videos with accompanying transcribed narrations. To achieve this goal, we propose a multilayer cross-modal attention network that enables effective optimization of a contrastive loss during training. We introduce a divided strategy that alternates between computing inter- and intra-modal attention across the visual and natural language modalities, which allows effective training via directly contrasting the two modalities' representations. We demonstrate the effectiveness of our approach by self-training on the HowTo100M instructional video dataset and evaluating on a newly collected dataset of localized described interactions in the YouCook2 dataset. We show that our approach outperforms alternative baselines, including shallow co-attention and full cross-modal attention. We also apply our approach to grounding phrases in images with weak supervision on Flickr30K and show that stacking multiple attention layers is effective and, when combined with a word-to-region loss, achieves state of the art on recall-at-one and pointing hand accuracies.

RETRIEVE: Coreset Selection for Efficient and Robust Semi-Supervised Learning

- Krishnateja Killamsetty, Xujiang Zhao, Feng Chen, Rishabh Iyer
- abstract@[open-review](#): Semi-supervised learning (SSL) algorithms have had great success in recent years in limited labeled data regimes. However, the current state-of-the-art SSL algorithms are computationally expensive and entail significant compute time and energy requirements. This can prove to be a huge limitation for many smaller companies and academic groups. Our main insight is that training on a subset of unlabeled data instead of entire unlabeled data enables the current SSL algorithms to converge faster, significantly reducing computational costs. In this work, we propose RETRIEVE, a coresnet selection framework for efficient and robust semi-supervised learning. RETRIEVE selects the coresnet by solving a mixed discrete-continuous bi-level optimization problem such that the selected coresnet minimizes the labeled set loss. We use a one-step gradient approximation and show that the discrete optimization problem is approximately submodular, enabling simple greedy algorithms to obtain the coresnet. We empirically demonstrate on several real-world datasets that existing SSL algorithms like VAT, Mean-Teacher, FixMatch, when used with RETRIEVE, achieve a) faster training times, b) better performance when unlabeled data consists of Out-of-Distribution (OOD) data and imbalance. More specifically, we show that with minimal accuracy degradation, RETRIEVE achieves a speedup of around \$3\times\$ in the traditional SSL setting and achieves a speedup of \$5\times\$ compared to state-of-the-art (SOTA) robust SSL algorithms in the case of imbalance and OOD data. RETRIEVE is available as a part of the CORDS toolkit: <https://github.com/decile-team/cords>.

Collaborating with Humans without Human Data

- DJ Strouse, Kevin McKee, Matt Botvinick, Edward Hughes, Richard Everett
- abstract@[open-review](#): Collaborating with humans requires rapidly adapting to their individual strengths, weaknesses, and preferences. Unfortunately, most standard multi-agent reinforcement learning techniques, such as self-play (SP) or population play (PP), produce agents that overfit to their training partners and do not generalize well to humans. Alternatively, researchers can collect human data, train a human model using behavioral cloning, and then use that model to train "human-aware" agents ("behavioral cloning play", or BCP). While such an approach can improve the generalization of agents to new human co-players, it involves the onerous and expensive step of collecting large amounts of human data first. Here, we study the problem of how to train agents that collaborate well with human partners without using human data. We argue that the crux of the problem is to produce a diverse set of training partners. Drawing inspiration from successful multi-agent approaches in competitive domains, we find that a surprisingly simple approach is

highly effective. We train our agent partner as the best response to a population of self-play agents and their past checkpoints taken throughout training, a method we call Fictitious Co-Play (FCP). Our experiments focus on a two-player collaborative cooking simulator that has recently been proposed as a challenge problem for coordination with humans. We find that FCP agents score significantly higher than SP, PP, and BCP when paired with novel agent and human partners. Furthermore, humans also report a strong subjective preference to partnering with FCP agents over all baselines.

[Training Feedback Spiking Neural Networks by Implicit Differentiation on the Equilibrium State](#)

- Mingqing Xiao, Qingyan Meng, Zongpeng Zhang, Yisen Wang, Zhouchen Lin
- abstract@[open-review](#): Spiking neural networks (SNNs) are brain-inspired models that enable energy-efficient implementation on neuromorphic hardware. However, the supervised training of SNNs remains a hard problem due to the discontinuity of the spiking neuron model. Most existing methods imitate the backpropagation framework and feedforward architectures for artificial neural networks, and use surrogate derivatives or compute gradients with respect to the spiking time to deal with the problem. These approaches either accumulate approximation errors or only propagate information limitedly through existing spikes, and usually require information propagation along time steps with large memory costs and biological implausibility. In this work, we consider feedback spiking neural networks, which are more brain-like, and propose a novel training method that does not rely on the exact reverse of the forward computation. First, we show that the average firing rates of SNNs with feedback connections would gradually evolve to an equilibrium state along time, which follows a fixed-point equation. Then by viewing the forward computation of feedback SNNs as a black-box solver for this equation, and leveraging the implicit differentiation on the equation, we can compute the gradient for parameters without considering the exact forward procedure. In this way, the forward and backward procedures are decoupled and therefore the problem of non-differentiable spiking functions is avoided. We also briefly discuss the biological plausibility of implicit differentiation, which only requires computing another equilibrium. Extensive experiments on MNIST, Fashion-MNIST, N-MNIST, CIFAR-10, and CIFAR-100 demonstrate the superior performance of our method for feedback models with fewer neurons and parameters in a small number of time steps. Our code is available at <https://github.com/pkuxmq/IDE-FSNN>.

[Online Selective Classification with Limited Feedback](#)

- Aditya Gangrade, Anil Kag, Ashok Cutkosky, Venkatesh Saligrama
- abstract@[open-review](#): Motivated by applications to resource-limited and safety-critical domains, we study selective classification in the online learning model, wherein a predictor may abstain from classifying an instance. For example, this may model an adaptive decision to invoke more resources on this instance. Two salient aspects of the setting we consider are that the data may be non-realisable, due to which abstention may be a valid long-term action, and that feedback is only received when the learner abstains, which models the fact that reliable labels are only available when the resource intensive processing is invoked. Within this framework, we explore strategies that make few mistakes, while not abstaining too many times more than the best-in-hindsight error-free classifier from a given class. That is, the one that makes no mistakes, while abstaining the fewest number of times. We construct simple versioning-based schemes for any $\mu \in (0,1]$ that make most T^μ mistakes while incurring $\tilde{O}(T^{1-\mu})$ excess abstention against adaptive adversaries. We further show that this dependence on T is tight, and provide illustrative experiments on realistic datasets.

[Controlled Text Generation as Continuous Optimization with Multiple Constraints](#)

- Sachin Kumar, Eric Malmi, Aliaksei Severyn, Yulia Tsvetkov
- abstract@[open-review](#): As large-scale language model pretraining pushes the state-of-the-art in text generation, recent work has turned to controlling attributes of the text such models generate. While modifying the pretrained models via fine-tuning remains the popular approach, it incurs a significant computational cost and can be infeasible due to a lack of appropriate data. As an alternative, we propose \textsc{MuCoCO}---a flexible and modular algorithm for controllable inference from pretrained models. We formulate the decoding process as an optimization problem that allows for multiple attributes we aim to control to be easily incorporated as differentiable constraints. By relaxing this discrete optimization to a continuous one, we make use of Lagrangian multipliers and gradient-descent-based techniques to generate the desired text. We evaluate our approach on controllable machine translation and style transfer with multiple sentence-level attributes and observe significant improvements over baselines.

[S\\$^3\\$: Sign-Sparse-Shift Reparametrization for Effective Training of Low-bit Shift Networks](#)

- Xinlin Li, Bang Liu, Yaoliang Yu, Wulong Liu, Chunjing XU, Vahid Partovi Nia
- abstract@[open-review](#): Shift neural networks reduce computation complexity by removing expensive multiplication operations and quantizing continuous weights into low-bit discrete values, which are fast and energy-efficient compared to conventional neural networks. However, existing shift networks are sensitive to the weight initialization and yield a degraded performance caused by vanishing gradient and weight sign freezing problem. To address these issues, we propose S\$^3\$ re-parameterization, a novel technique for training low-bit shift networks. Our method decomposes a discrete parameter in a sign-sparse-shift 3-fold manner. This way, it efficiently learns a low-bit network with weight dynamics similar to full-precision networks and insensitive to weight initialization. Our proposed training method pushes the boundaries of shift neural networks and shows 3-bit shift networks compete with their full-precision counterparts in terms of top-1 accuracy on ImageNet.

[Implicit MLE: Backpropagating Through Discrete Exponential Family Distributions](#)

- Mathias Niepert, Pasquale Minervini, Luca Franceschi
- abstract@[open-review](#): Combining discrete probability distributions and combinatorial optimization problems with neural network components has numerous applications but poses several challenges. We propose Implicit Maximum Likelihood Estimation (I-MLE), a framework for end-to-end learning of models combining discrete exponential family distributions and differentiable neural components. I-MLE is widely applicable as it only requires the ability to compute the most probable states and does not rely on smooth relaxations. The framework encompasses several approaches such as perturbation-based implicit differentiation and recent methods to differentiate through black-box combinatorial solvers. We introduce a novel class of noise distributions for approximating marginals via perturb-and-MAP. Moreover, we show that I-MLE simplifies to maximum likelihood estimation when used in some recently studied learning settings that involve combinatorial solvers. Experiments on several datasets suggest that I-MLE is competitive with and often outperforms existing approaches which rely on problem-specific relaxations.

[Scaling up Continuous-Time Markov Chains Helps Resolve Underspecification](#)

- Alkis Gotovos, Rebekka Burkholz, John Quackenbush, Stefanie Jegelka
- abstract@[open-review](#): Modeling the time evolution of discrete sets of items (e.g., genetic mutations) is a fundamental problem in many biomedical applications. We approach this problem through the lens of continuous-time Markov chains, and show that the resulting learning task is generally underspecified in the usual setting of cross-sectional data. We explore a perhaps surprising remedy: including a number of additional independent items can help determine time order, and hence resolve underspecification. This is in sharp contrast to the common practice of limiting the analysis to a small subset of relevant items, which is followed largely due to poor scaling of existing methods. To put our theoretical insight into practice, we develop an approximate likelihood maximization method for learning continuous-time Markov chains, which can scale to hundreds of items and is orders of magnitude faster than previous methods. We demonstrate the effectiveness of our approach on synthetic and real cancer data.

[Do Neural Optimal Transport Solvers Work? A Continuous Wasserstein-2 Benchmark](#)

- Alexander Korotin, Lingxiao Li, Aude Genevay, Justin M. Solomon, Alexander Filippov, Evgeny Burnaev
- abstract@[open-review](#): Despite the recent popularity of neural network-based solvers for optimal transport (OT), there is no standard quantitative way to evaluate their performance. In this paper, we address this issue for quadratic-cost transport---specifically, computation of the Wasserstein-2 distance, a commonly-used formulation of optimal transport in machine learning. To overcome the challenge of computing ground truth transport maps between continuous measures needed to assess these solvers, we use input-convex neural networks (ICNN) to construct pairs of measures whose ground truth OT maps can be obtained analytically. This strategy yields pairs of continuous benchmark measures in high-dimensional spaces such as spaces of images. We thoroughly evaluate existing optimal transport solvers using these benchmark measures. Even though these solvers perform well in downstream tasks, many do not faithfully recover optimal transport maps. To investigate the cause of this discrepancy, we further test the solvers in a setting of image generation. Our study reveals crucial limitations of existing solvers and shows that increased OT accuracy does not necessarily correlate to better results downstream.

[Linear Convergence in Federated Learning: Tackling Client Heterogeneity and Sparse Gradients](#)

- Aritra Mitra, Rayana Jaafar, George J. Pappas, Hamed Hassani
- abstract@[open-review](#): We consider a standard federated learning (FL) setup where a group of clients periodically coordinate with a central server to train a statistical model. We develop a general algorithmic framework called FedLin to tackle some of the key challenges intrinsic to FL, namely objective heterogeneity, systems heterogeneity, and infrequent and imprecise communication. Our framework is motivated by the observation that under these challenges, various existing FL algorithms suffer from a fundamental speed-accuracy conflict: they either guarantee linear convergence but to an incorrect point, or convergence to the global minimum but at a sub-linear rate, i.e., fast convergence comes at the expense of accuracy. In contrast, when the clients' local loss functions are smooth and strongly convex, we show that FedLin guarantees linear convergence to the global minimum, despite arbitrary objective and systems heterogeneity. We then establish matching upper and lower bounds on the convergence rate of FedLin that highlight the effects of infrequent, periodic communication. Finally, we show that FedLin preserves linear convergence rates under aggressive gradient sparsification, and quantify the effect of the compression level on the convergence rate. Notably, our work is the first to provide tight linear convergence rate guarantees, and constitutes the first comprehensive analysis of gradient sparsification in FL.

[On the Convergence of Prior-Guided Zeroth-Order Optimization Algorithms](#)

- Shuyu Cheng, Guoqiang Wu, Jun Zhu
- abstract@[open-review](#): Zeroth-order (ZO) optimization is widely used to handle challenging tasks, such as query-based black-box adversarial attacks and reinforcement learning. Various attempts have been made to integrate prior information into the gradient estimation procedure based on finite differences, with promising empirical results. However, their convergence properties are not well understood. This paper makes an attempt to fill up this gap by analyzing the convergence of prior-guided ZO algorithms under a greedy descent framework with various gradient estimators. We provide a convergence guarantee for the prior-guided random gradient-free (PRGF) algorithms. Moreover, to further accelerate over greedy descent methods, we present a new accelerated random search (ARS) algorithm that incorporates prior information, together with a convergence analysis. Finally, our theoretical results are confirmed by experiments on several numerical benchmarks as well as adversarial attacks.

[Revisit Multimodal Meta-Learning through the Lens of Multi-Task Learning](#)

- Milad Abdollahzadeh, Touba Malekzadeh, Ngai-Man (Man) Cheung
- abstract@[open-review](#): Multimodal meta-learning is a recent problem that extends conventional few-shot meta-learning by generalizing its setup to diverse multimodal task distributions. This setup makes a step towards mimicking how humans make use of a diverse set of prior skills to learn new skills. Previous work has achieved encouraging performance. In particular, in spite of the diversity of the multimodal tasks, previous work claims that a single meta-learner trained on a multimodal distribution can sometimes outperform multiple specialized meta-learners trained on individual unimodal distributions. The improvement is attributed to knowledge transfer between different modes of task distributions. However, there is no deep investigation to verify and understand the knowledge transfer between multimodal tasks. Our work makes two contributions to multimodal meta-learning. First, we propose a method to quantify knowledge transfer between tasks of different modes at a micro-level. Our quantitative, task-level analysis is inspired by the recent transference idea from multi-task learning. Second, inspired by hard parameter sharing in multi-task learning and a new interpretation of related work, we propose a new multimodal meta-learner that outperforms existing work by considerable margins. While the major focus is on multimodal meta-learning, our work also attempts to shed light on task interaction in conventional meta-learning. The code for this project is available at <https://miladabd.github.io/KML>.

[Dynamic Sasvi: Strong Safe Screening for Norm-Regularized Least Squares](#)

- Hiroaki Yamada, Makoto Yamada
- abstract@[open-review](#): A recently introduced technique, called safe screening," for a sparse optimization problem allows us to identify irrelevant variables in the early stages of optimization. In this paper, we first propose a flexible framework for safe screening based on the Fenchel--Rockafellar duality and then derive a strong safe screening rule for norm-regularized least squares using the proposed framework. We refer to the proposed screening rule for norm-regularized least squares as "dynamic Sasvi" because it can be interpreted as a generalization of Sasvi. Unlike the original Sasvi, it does not require the exact solution of a more strongly regularized problem; hence, it works safely in practice. We show that our screening rule always eliminates more features compared with the existing state-of-the-art methods.

[What Matters for Adversarial Imitation Learning?](#)

- Manu Orsini, Anton Raichuk, Leonard Huszenot, Damien Vincent, Robert Dadashi, Sertan Girgin, Matthieu Geist, Olivier Bachem, Olivier Pietquin, Marcin Andrychowicz
- abstract@[open-review](#): Adversarial imitation learning has become a popular framework for imitation in continuous control. Over the years, several variations of its components were proposed to enhance the performance of the learned policies as well as the sample complexity of the algorithm. In practice, these choices are rarely tested all together in rigorous empirical studies. It is therefore difficult to discuss and understand what choices, among the high-level algorithmic options as well as low-level implementation details, matter. To tackle this issue, we implement more than 50 of these choices in a generic adversarial imitation learning framework and investigate their impacts in a large-scale study (>500k trained agents) with both synthetic and human-generated demonstrations. We analyze the key results and highlight the most surprising findings.

[Sequential Causal Imitation Learning with Unobserved Confounders](#)

- Daniel Kumor, Junzhe Zhang, Elias Bareinboim
- abstract@[open-review](#): "Monkey see monkey do" is an age-old adage, referring to naive imitation without a deep understanding of a system's underlying mechanics. Indeed, if a demonstrator has access to information unavailable to the imitator (monkey), such as a different set of sensors, then no matter how perfectly the imitator models its perceived environment (See), attempting to directly reproduce the demonstrator's behavior (Do) can lead to poor outcomes. Imitation learning in the presence of a mismatch between demonstrator and imitator has been studied in the literature under the rubric of causal imitation learning (Zhang et. al. 2020), but existing solutions are limited to single-stage decision-making. This paper investigates the problem of causal imitation learning in sequential settings, where the imitator must make multiple decisions per episode. We develop a graphical criterion that is both

necessary and sufficient for determining the feasibility of causal imitation, providing conditions when an imitator can match a demonstrator's performance despite differing capabilities. Finally, we provide an efficient algorithm for determining imitability, and corroborate our theory with simulations.

[Topic Modeling Revisited: A Document Graph-based Neural Network Perspective](#)

- Dazhong Shen, Chuan Qin, Chao Wang, Zheng Dong, Hengshu Zhu, Hui Xiong
- abstract@[open-review](#): Most topic modeling approaches are based on the bag-of-words assumption, where each word is required to be conditionally independent in the same document. As a result, both of the generative story and the topic formulation have totally ignored the semantic dependency among words, which is important for improving the semantic comprehension and model interpretability. To this end, in this paper, we revisit the task of topic modeling by transforming each document into a directed graph with word dependency as edges between word nodes, and develop a novel approach, namely Graph Neural Topic Model (GNTM). Specifically, in GNTM, a well-defined probabilistic generative story is designed to model both the graph structure and word sets with multinomial distributions on the vocabulary and word dependency edge set as the topics. Meanwhile, a Neural Variational Inference (NVI) approach is proposed to learn our model with graph neural networks to encode the document graphs. Besides, we theoretically demonstrate that Latent Dirichlet Allocation (LDA) can be derived from GNTM as a special case with similar objective functions. Finally, extensive experiments on four benchmark datasets have clearly demonstrated the effectiveness and interpretability of GNTM compared with state-of-the-art baselines.

[Hard-Attention for Scalable Image Classification](#)

- Athanasios Papadopoulos, Paweł Korus, Nasir Memon
- abstract@[open-review](#): Can we leverage high-resolution information without the unsustainable quadratic complexity to input scale? We propose Traversal Network (TNet), a novel multi-scale hard-attention architecture, which traverses image scale-space in a top-down fashion, visiting only the most informative image regions along the way. TNet offers an adjustable trade-off between accuracy and complexity, by changing the number of attended image locations. We compare our model against hard-attention baselines on ImageNet, achieving higher accuracy with less resources (FLOPs, processing time and memory). We further test our model on fMoW dataset, where we process satellite images of size up to \$896 \times 896\$ px, getting up to \$2.5\times\$ faster processing compared to baselines operating on the same resolution, while achieving higher accuracy as well. TNet is modular, meaning that most classification models could be adopted as its backbone for feature extraction, making the reported performance gains orthogonal to benefits offered by existing optimized deep models. Finally, hard-attention guarantees a degree of interpretability to our model's predictions, without any extra cost beyond inference.

[Fast Routing under Uncertainty: Adaptive Learning in Congestion Games via Exponential Weights](#)

- Dong Quan Vu, Kimon Antonakopoulos, Panayotis Mertikopoulos
- abstract@[open-review](#): We examine an adaptive learning framework for nonatomic congestion games where the players' cost functions may be subject to exogenous fluctuations (e.g., due to disturbances in the network, variations in the traffic going through a link). In this setting, the popular multiplicative/exponential weights algorithm enjoys an $\mathcal{O}(1/\sqrt{T})$ equilibrium convergence rate; however, this rate is suboptimal in static environments---i.e., when the network is not subject to randomness. In this static regime, accelerated algorithms achieve an $\mathcal{O}(1/T^{1/2})$ convergence speed, but they fail to converge altogether in stochastic problems. To fill this gap, we propose a novel, adaptive exponential weights method---dubbed AdaWeight---that seamlessly interpolates between the $\mathcal{O}(1/T^{1/2})$ and $\mathcal{O}(1/\sqrt{T})$ rates in the static and stochastic regimes respectively. Importantly, this "best-of-both-worlds" guarantee does not require any prior knowledge of the problem's parameters or tuning by the optimizer; in addition, the method's convergence speed depends subquadratically on the size of the network (number of vertices and edges), so it scales gracefully to large, real-life urban networks.

[Profiling Pareto Front With Multi-Objective Stein Variational Gradient Descent](#)

- Xingchao Liu, Xin Tong, Qiang Liu
- abstract@[open-review](#): Finding diverse and representative Pareto solutions from the Pareto front is a key challenge in multi-objective optimization (MOO). In this work, we propose a novel gradient-based algorithm for profiling Pareto front by using Stein variational gradient descent (SVGD). We also provide a counterpart of our method based on Langevin dynamics. Our methods iteratively update a set of points in a parallel fashion to push them towards the Pareto front using multiple gradient descent, while encouraging the diversity between the particles by using the repulsive force mechanism in SVGD, or diffusion noise in Langevin dynamics. Compared with existing gradient-based methods that require predefined preference functions, our method can work efficiently in high dimensional problems, and can obtain more diverse solutions evenly distributed in the Pareto front. Moreover, our methods are theoretically guaranteed to converge to the Pareto front. We demonstrate the effectiveness of our method, especially the SVGD algorithm, through extensive experiments, showing its superiority over existing gradient-based algorithms.

[MAP Propagation Algorithm: Faster Learning with a Team of Reinforcement Learning Agents](#)

- Stephen Chung
- abstract@[open-review](#): Nearly all state-of-the-art deep learning algorithms rely on error backpropagation, which is generally regarded as biologically implausible. An alternative way of training an artificial neural network is through treating each unit in the network as a reinforcement learning agent, and thus the network is considered as a team of agents. As such, all units can be trained by REINFORCE, a local learning rule modulated by a global signal that is more consistent with biologically observed forms of synaptic plasticity. Although this learning rule follows the gradient of return in expectation, it suffers from high variance and thus the low speed of learning, rendering it impractical to train deep networks. We therefore propose a novel algorithm called MAP propagation to reduce this variance significantly while retaining the local property of the learning rule. Experiments demonstrated that MAP propagation could solve common reinforcement learning tasks at a similar speed to backpropagation when applied to an actor-critic network. Our work thus allows for the broader application of teams of agents in deep reinforcement learning.

[TransGAN: Two Pure Transformers Can Make One Strong GAN, and That Can Scale Up](#)

- Yifan Jiang, Shiyu Chang, Zhangyang Wang
- abstract@[open-review](#): The recent explosive interest on transformers has suggested their potential to become powerful ``universal'' models for computer vision tasks, such as classification, detection, and segmentation. While those attempts mainly study the discriminative models, we explore transformers on some more notoriously difficult vision tasks, e.g., generative adversarial networks (GANs). Our goal is to conduct the first pilot study in building a GAN completely free of convolutions, using only pure transformer-based architectures. Our vanilla GAN architecture, dubbed TransGAN, consists of a memory-friendly transformer-based generator that progressively increases feature resolution, and correspondingly a multi-scale discriminator to capture simultaneously semantic contexts and low-level textures. On top of them, we introduce the new module of grid self-attention for alleviating the memory bottleneck further, in order to scale up TransGAN to high-resolution generation. We also develop a unique training recipe including a series of techniques that can mitigate the training instability issues of TransGAN, such as data augmentation, modified normalization, and relative position encoding. Our best architecture achieves highly competitive performance compared to current state-of-the-art GANs using convolutional backbones. Specifically, TransGAN sets new state-of-the-art inception score of 10.43 and FID of 18.28 on STL-10. It also reaches the inception score of 9.02 and FID of 9.26 on CIFAR-10, and 5.28 FID on CelebA 128×128 , respectively: both on par with the current best results and outperforming StyleGAN-V2. When it comes to higher-resolution (e.g. 256×256) generation tasks, such as on CelebA-HQ

and LSUN-Church, TransGAN continues to produce diverse visual examples with high fidelity and impressive texture details. In addition, we dive deep into the transformer-based generation models to understand how their behaviors differ from convolutional ones, by visualizing training dynamics. The code is available at: <https://github.com/VITA-Group/TransGAN>.

[A Central Limit Theorem for Differentially Private Query Answering](#)

- Jinshuo Dong, Weijie Su, Linjun Zhang
- abstract@[open-review](#): Perhaps the single most important use case for differential privacy is to privately answer numerical queries, which is usually achieved by adding noise to the answer vector. The central question is, therefore, to understand which noise distribution optimizes the privacy-accuracy trade-off, especially when the dimension of the answer vector is high. Accordingly, an extensive literature has been dedicated to the question and the upper and lower bounds have been successfully matched up to constant factors (Bun et al., 2018; Steinke & Ullman, 2017). In this paper, we take a novel approach to address this important optimality question. We first demonstrate an intriguing central limit theorem phenomenon in the high-dimensional regime. More precisely, we prove that a mechanism is approximately Gaussian Differentially Private (Dong et al., 2021) if the added noise satisfies certain conditions. In particular, densities proportional to $\mathrm{e}^{-\|\mathbf{x}\|_p^\alpha}$, where $\|\mathbf{x}\|_p$ is the standard ℓ_p -norm, satisfies the conditions. Taking this perspective, we make use of the Cramer--Rao inequality and show an "uncertainty principle"-style result: the product of privacy parameter and the ℓ_2 -loss of the mechanism is lower bounded by the dimension. Furthermore, the Gaussian mechanism achieves the constant-sharp optimal privacy-accuracy trade-off among all such noises. Our findings are corroborated by numerical experiments.

[Differential Privacy Dynamics of Langevin Diffusion and Noisy Gradient Descent](#)

- Rishav Chourasia, Jiayuan Ye, Reza Shokri
- abstract@[open-review](#): What is the information leakage of an iterative randomized learning algorithm about its training data, when the internal state of the algorithm is \emph{private}? How much is the contribution of each specific training epoch to the information leakage through the released model? We study this problem for noisy gradient descent algorithms, and model the \emph{dynamics} of R\enyi differential privacy loss throughout the training process. Our analysis traces a provably \emph{tight} bound on the R\enyi divergence between the pair of probability distributions over parameters of models trained on neighboring datasets. We prove that the privacy loss converges exponentially fast, for smooth and strongly convex loss functions, which is a significant improvement over composition theorems (which over-estimate the privacy loss by upper-bounding its total value over all intermediate gradient computations). For Lipschitz, smooth, and strongly convex loss functions, we prove optimal utility with a small gradient complexity for noisy gradient descent algorithms.

[Data driven semi-supervised learning](#)

- Maria-Florina F. Balcan, Dravyansh Sharma
- abstract@[open-review](#): We consider a novel data driven approach for designing semi-supervised learning algorithms that can effectively learn with only a small number of labeled examples. We focus on graph-based techniques, where the unlabeled examples are connected in a graph under the implicit assumption that similar nodes likely have similar labels. Over the past two decades, several elegant graph-based semi-supervised learning algorithms for inferring the labels of the unlabeled examples given the graph and a few labeled examples have been proposed. However, the problem of how to create the graph (which impacts the practical usefulness of these methods significantly) has been relegated to heuristics and domain-specific art, and no general principles have been proposed. In this work we present a novel data driven approach for learning the graph and provide strong formal guarantees in both the distributional and online learning formalizations. We show how to leverage problem instances coming from an underlying problem domain to learn the graph hyperparameters for commonly used parametric families of graphs that provably perform well on new instances from the same domain. We obtain low regret and efficient algorithms in the online setting, and generalization guarantees in the distributional setting. We also show how to combine several very different similarity metrics and learn multiple hyperparameters, our results hold for large classes of problems. We expect some of the tools and techniques we develop along the way to be of independent interest, for data driven algorithms more generally.

[Online Meta-Learning via Learning with Layer-Distributed Memory](#)

- Sudarshan Babu, Pedro Savarese, Michael Maire
- abstract@[open-review](#): We demonstrate that efficient meta-learning can be achieved via end-to-end training of deep neural networks with memory distributed across layers. The persistent state of this memory assumes the entire burden of guiding task adaptation. Moreover, its distributed nature is instrumental in orchestrating adaptation. Ablation experiments demonstrate that providing relevant feedback to memory units distributed across the depth of the network enables them to guide adaptation throughout the entire network. Our results show that this is a successful strategy for simplifying meta-learning -- often cast as a bi-level optimization problem -- to standard end-to-end training, while outperforming gradient-based, prototype-based, and other memory-based meta-learning strategies. Additionally, our adaptation strategy naturally handles online learning scenarios with a significant delay between observing a sample and its corresponding label -- a setting in which other approaches struggle. Adaptation via distributed memory is effective across a wide range of learning tasks, ranging from classification to online few-shot semantic segmentation.

[Physics-Integrated Variational Autoencoders for Robust and Interpretable Generative Modeling](#)

- Naoya Takeishi, Alexandros Kalousis
- abstract@[open-review](#): Integrating physics models within machine learning models holds considerable promise toward learning robust models with improved interpretability and abilities to extrapolate. In this work, we focus on the integration of incomplete physics models into deep generative models. In particular, we introduce an architecture of variational autoencoders (VAEs) in which a part of the latent space is grounded by physics. A key technical challenge is to strike a balance between the incomplete physics and trainable components such as neural networks for ensuring that the physics part is used in a meaningful manner. To this end, we propose a regularized learning method that controls the effect of the trainable components and preserves the semantics of the physics-based latent variables as intended. We not only demonstrate generative performance improvements over a set of synthetic and real-world datasets, but we also show that we learn robust models that can consistently extrapolate beyond the training distribution in a meaningful manner. Moreover, we show that we can control the generative process in an interpretable manner.

[Characterizing the risk of fairwashing](#)

- Ulrich Aïvodji, Hiromi Arai, Sébastien Gambs, Satoshi Hara
- abstract@[open-review](#): Fairwashing refers to the risk that an unfair black-box model can be explained by a fairer model through post-hoc explanation manipulation. In this paper, we investigate the capability of fairwashing attacks by analyzing their fidelity-unfairness trade-offs. In particular, we show that fairwashed explanation models can generalize beyond the suing group (i.e., data points that are being explained), meaning that a fairwashed explainer can be used to rationalize subsequent unfair decisions of a black-box model. We also demonstrate that fairwashing attacks can transfer across black-box models, meaning that other black-box models can perform fairwashing without explicitly using their predictions. This generalization and transferability of fairwashing attacks imply that their detection will be difficult in practice. Finally, we propose an approach to quantify the risk of fairwashing, which is based on the computation of the range of the unfairness of high-fidelity explainers.

[Qimera: Data-free Quantization with Synthetic Boundary Supporting Samples](#)

- Kanghyun Choi, Deokki Hong, Noseong Park, Youngsok Kim, Jinho Lee
- abstract@[open-review](#): Model quantization is known as a promising method to compress deep neural networks, especially for inferences on lightweight mobile or edge devices. However, model quantization usually requires access to the original training data to maintain the accuracy of the full-precision models, which is often infeasible in real-world scenarios for security and privacy issues. A popular approach to perform quantization without access to the original data is to use synthetically generated samples, based on batch-normalization statistics or adversarial learning. However, the drawback of such approaches is that they primarily rely on random noise input to the generator to attain diversity of the synthetic samples. We find that this is often insufficient to capture the distribution of the original data, especially around the decision boundaries. To this end, we propose Qimera, a method that uses superposed latent embeddings to generate synthetic boundary supporting samples. For the superposed embeddings to better reflect the original distribution, we also propose using an additional disentanglement mapping layer and extracting information from the full-precision model. The experimental results show that Qimera achieves state-of-the-art performances for various settings on data-free quantization. Code is available at <https://github.com/iamkanghyunchoi/qimera>.

[Embedding Principle of Loss Landscape of Deep Neural Networks](#)

- Yaoyu Zhang, Zhongwang Zhang, Tao Luo, Zhiqin J Xu
- abstract@[open-review](#): Understanding the structure of loss landscape of deep neural networks (DNNs) is obviously important. In this work, we prove an embedding principle that the loss landscape of a DNN "contains" all the critical points of all the narrower DNNs. More precisely, we propose a critical embedding such that any critical point, e.g., local or global minima, of a narrower DNN can be embedded to a critical point/affine subspace of the target DNN with higher degeneracy and preserving the DNN output function. Note that, given any training data, differentiable loss function and differentiable activation function, this embedding structure of critical points holds. This general structure of DNNs is starkly different from other nonconvex problems such as protein-folding. Empirically, we find that a wide DNN is often attracted by highly-degenerate critical points that are embedded from narrow DNNs. The embedding principle provides a new perspective to study the general easy optimization of wide DNNs and unravels a potential implicit low-complexity regularization during the training. Overall, our work provides a skeleton for the study of loss landscape of DNNs and its implication, by which a more exact and comprehensive understanding can be anticipated in the near future.

[Adversarial Reweighting for Partial Domain Adaptation](#)

- Xiang Gu, Xi Yu, yan yang, Jian Sun, Zongben Xu
- abstract@[open-review](#): Partial domain adaptation (PDA) has gained much attention due to its practical setting. The current PDA methods usually adapt the feature extractor by aligning the target and reweighted source domain distributions. In this paper, we experimentally find that the feature adaptation by the reweighted distribution alignment in some state-of-the-art PDA methods is not robust to the ``noisy'' weights of source domain data, leading to negative domain transfer on some challenging benchmarks. To tackle the challenge of negative domain transfer, we propose a novel Adversarial Reweighting (AR) approach that adversarially learns the weights of source domain data to align the source and target domain distributions, and the transferable deep recognition network is learned on the reweighted source domain data. Based on this idea, we propose a training algorithm that alternately updates the parameters of the network and optimizes the weights of source domain data. Extensive experiments show that our method achieves state-of-the-art results on the benchmarks of ImageNet-Caltech, Office-Home, VisDA-2017, and DomainNet. Ablation studies also confirm the effectiveness of our approach.

[M-FAC: Efficient Matrix-Free Approximations of Second-Order Information](#)

- Elias Frantar, Eldar Kurtic, Dan Alistarh
- abstract@[open-review](#): Efficiently approximating local curvature information of the loss function is a useful tool for the optimization and compression of deep neural networks. Yet, most existing methods to approximate second-order information have high computational or storage costs, limiting their practicality. In this work, we investigate matrix-free approaches for estimating Inverse-Hessian Vector Products (IHVPs) for the case when the Hessian can be approximated as a sum of rank-one matrices, as in the classic approximation of the Hessian by the empirical Fisher matrix. The first algorithm we propose is tailored towards network compression and can compute the IHVP for dimension d given a fixed set of m rank-one matrices using $\mathcal{O}(dm^2)$ precomputation, $\mathcal{O}(dm)$ cost for computing the IHVP and query cost $\mathcal{O}(m)$ for computing any single element of the inverse Hessian approximation. The second algorithm targets an optimization setting, where we wish to compute the product between the inverse Hessian, estimated over a sliding window of optimization steps, and a given gradient direction. We give an algorithm with cost $\mathcal{O}(dm + m^2)$ for computing the IHVP and $\mathcal{O}(dm + m^3)$ for adding or removing any gradient from the sliding window. We show that both algorithms yield competitive results for network pruning and optimization, respectively, with significantly lower computational overhead relative to existing second-order methods.

[Graph Adversarial Self-Supervised Learning](#)

- Longqi Yang, Liangliang Zhang, Wenjing Yang
- abstract@[open-review](#): This paper studies a long-standing problem of learning the representations of a whole graph without human supervision. The recent self-supervised learning methods train models to be invariant to the transformations (views) of the inputs. However, designing these views requires the experience of human experts. Inspired by adversarial training, we propose an adversarial self-supervised learning (`GASSL`) framework for learning unsupervised representations of graph data without any handcrafted views. `GASSL` automatically generates challenging views by adding perturbations to the input and are adversarially trained with respect to the encoder. Our method optimizes the min-max problem and utilizes a gradient accumulation strategy to accelerate the training process. Experimental on ten graph classification datasets show that the proposed approach is superior to state-of-the-art self-supervised learning baselines, which are competitive with supervised models.

[Anti-Backdoor Learning: Training Clean Models on Poisoned Data](#)

- Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, Xingjun Ma
- abstract@[open-review](#): Backdoor attack has emerged as a major security threat to deep neural networks (DNNs). While existing defense methods have demonstrated promising results on detecting or erasing backdoors, it is still not clear whether robust training methods can be devised to prevent the backdoor triggers being injected into the trained model in the first place. In this paper, we introduce the concept of `anti-backdoor learning`, aiming to train `clean` models given backdoor-poisoned data. We frame the overall learning process as a dual-task of learning the `clean` and the `backdoor` portions of data. From this view, we identify two inherent characteristics of backdoor attacks as their weaknesses: 1) the models learn backdoored data much faster than learning with clean data, and the stronger the attack the faster the model converges on backdoored data; 2) the backdoor task is tied to a specific class (the backdoor target class). Based on these two weaknesses, we propose a general learning scheme, Anti-Backdoor Learning (ABL), to automatically prevent backdoor attacks during training. ABL introduces a two-stage `gradient ascent` mechanism for standard training to 1) help isolate backdoor examples at an early training stage, and 2) break the correlation between backdoor examples and the target class at a later training stage. Through extensive experiments on multiple benchmark datasets against 10 state-of-the-art attacks, we empirically show that ABL-trained models on backdoor-poisoned data achieve the same performance as they were trained on purely clean data. Code is available at <https://github.com/bboylg/ABL>.

[Locally Most Powerful Bayesian Test for Out-of-Distribution Detection using Deep Generative Models](#)

- Keunseo Kim, JunCheol Shin, Heeyoung Kim

- abstract@[open-review](#): Several out-of-distribution (OOD) detection scores have been recently proposed for deep generative models because the direct use of the likelihood threshold for OOD detection has been shown to be problematic. In this paper, we propose a new OOD score based on a Bayesian hypothesis test called the locally most powerful Bayesian test (LMPBT). The LMPBT is locally most powerful in that the alternative hypothesis (the representative parameter for the OOD sample) is specified to maximize the probability that the Bayes factor exceeds the evidence threshold in favor of the alternative hypothesis provided that the parameter specified under the alternative hypothesis is in the neighborhood of the parameter specified under the null hypothesis. That is, under this neighborhood parameter condition, the test with the proposed alternative hypothesis maximizes the probability of correct detection of OOD samples. We also propose numerical strategies for more efficient and reliable computation of the LMPBT for practical application to deep generative models. Evaluations conducted of the OOD detection performance of the LMPBT on various benchmark datasets demonstrate its superior performance over existing OOD detection methods.

[Stable Neural ODE with Lyapunov-Stable Equilibrium Points for Defending Against Adversarial Attacks](#)

- Qiyu Kang, Yang Song, Qinxu Ding, Wee Peng Tay
- abstract@[open-review](#): Deep neural networks (DNNs) are well-known to be vulnerable to adversarial attacks, where malicious human-imperceptible perturbations are included in the input to the deep network to fool it into making a wrong classification. Recent studies have demonstrated that neural Ordinary Differential Equations (ODEs) are intrinsically more robust against adversarial attacks compared to vanilla DNNs. In this work, we propose a neural ODE with Lyapunov-stable equilibrium points for defending against adversarial attacks (SODEF). By ensuring that the equilibrium points of the ODE solution used as part of SODEF are Lyapunov-stable, the ODE solution for an input with a small perturbation converges to the same solution as the unperturbed input. We provide theoretical results that give insights into the stability of SODEF as well as the choice of regularizers to ensure its stability. Our analysis suggests that our proposed regularizers force the extracted feature points to be within a neighborhood of the Lyapunov-stable equilibrium points of the SODEF ODE. SODEF is compatible with many defense methods and can be applied to any neural network's final regressor layer to enhance its stability against adversarial attacks.

[Robust Compressed Sensing MRI with Deep Generative Priors](#)

- Ajil Jalal, Marius Arvinte, Giannis Daras, Eric Price, Alexandros G. Dimakis, Jon Tamir
- abstract@[open-review](#): The CSGM framework (Bora-Jalal-Price-Dimakis'17) has shown that deep generative priors can be powerful tools for solving inverse problems. However, to date this framework has been empirically successful only on certain datasets (for example, human faces and MNIST digits), and it is known to perform poorly on out-of-distribution samples. In this paper, we present the first successful application of the CSGM framework on clinical MRI data. We train a generative prior on brainscans from the fastMRI dataset, and show that posterior sampling via Langevin dynamics achieves high quality reconstructions. Furthermore, our experiments and theory show that posterior sampling is robust to changes in the ground-truth distribution and measurement process. Our code and models are available at: \url{https://github.com/utcsilab/csgm-mri-langevin}.

[H-NeRF: Neural Radiance Fields for Rendering and Temporal Reconstruction of Humans in Motion](#)

- Hongyi Xu, Thiem Alldieck, Cristian Sminchisescu
- abstract@[open-review](#): We present neural radiance fields for rendering and temporal (4D) reconstruction of humans in motion (H-NeRF), as captured by a sparse set of cameras or even from a monocular video. Our approach combines ideas from neural scene representation, novel-view synthesis, and implicit statistical geometric human representations, coupled using novel loss functions. Instead of learning a radiance field with a uniform occupancy prior, we constrain it by a structured implicit human body model, represented using signed distance functions. This allows us to robustly fuse information from sparse views and generalize well beyond the poses or views observed in training. Moreover, we apply geometric constraints to co-learn the structure of the observed subject -- including both body and clothing -- and to regularize the radiance field to geometrically plausible solutions. Extensive experiments on multiple datasets demonstrate the robustness and the accuracy of our approach, its generalization capabilities significantly outside a small training set of poses and views, and statistical extrapolation beyond the observed shape.

[DOBF: A Deobfuscation Pre-Training Objective for Programming Languages](#)

- Marie-Anne Lachaux, Baptiste Roziere, Marc Szafraniec, Guillaume Lample
- abstract@[open-review](#): Recent advances in self-supervised learning have dramatically improved the state of the art on a wide variety of tasks. However, research in language model pre-training has mostly focused on natural languages, and it is unclear whether models like BERT and its variants provide the best pre-training when applied to other modalities, such as source code. In this paper, we introduce a new pre-training objective, DOBF, that leverages the structural aspect of programming languages and pre-trains a model to recover the original version of obfuscated source code. We show that models pre-trained with DOBF significantly outperform existing approaches on multiple downstream tasks, providing relative improvements of up to 12.2% in unsupervised code translation, and 5.3% in natural language code search. Incidentally, we found that our pre-trained model is able to deobfuscate fully obfuscated source files, and to suggest descriptive variable names.

[Detecting Errors and Estimating Accuracy on Unlabeled Data with Self-training Ensembles](#)

- Jiefeng Chen, Frederick Liu, Besim Avci, Xi Wu, Yingyu Liang, Somesh Jha
- abstract@[open-review](#): When a deep learning model is deployed in the wild, it can encounter test data drawn from distributions different from the training data distribution and suffer drop in performance. For safe deployment, it is essential to estimate the accuracy of the pre-trained model on the test data. However, the labels for the test inputs are usually not immediately available in practice, and obtaining them can be expensive. This observation leads to two challenging tasks: (1) unsupervised accuracy estimation, which aims to estimate the accuracy of a pre-trained classifier on a set of unlabeled test inputs; (2) error detection, which aims to identify mis-classified test inputs. In this paper, we propose a principled and practically effective framework that simultaneously addresses the two tasks. The proposed framework iteratively learns an ensemble of models to identify mis-classified data points and performs self-training to improve the ensemble with the identified points. Theoretical analysis demonstrates that our framework enjoys provable guarantees for both accuracy estimation and error detection under mild conditions readily satisfied by practical deep learning models. Along with the framework, we proposed and experimented with two instantiations and achieved state-of-the-art results on 59 tasks. For example, on iWildCam, one instantiation reduces the estimation error for unsupervised accuracy estimation by at least 70% and improves the F1 score for error detection by at least 4.7% compared to existing methods.

[Exploiting Chain Rule and Bayes' Theorem to Compare Probability Distributions](#)

- Huangjie Zheng, Mingyuan Zhou
- abstract@[open-review](#): To measure the difference between two probability distributions, referred to as the source and target, respectively, we exploit both the chain rule and Bayes' theorem to construct conditional transport (CT), which is constituted by both a forward component and a backward one. The forward CT is the expected cost of moving a source data point to a target one, with their joint distribution defined by the product of the source probability density function (PDF) and a source-dependent conditional distribution, which is related to the target PDF via Bayes' theorem. The backward CT is defined by reversing the direction. The CT cost can be approximated by replacing the source and target PDFs with their discrete empirical distributions supported on mini-batches, making it amenable to implicit distributions and stochastic gradient descent-based optimization. When applied to train a generative model, CT is shown to strike a good balance between mode-covering and mode-seeking behaviors and strongly resist mode collapse. On a wide

variety of benchmark datasets for generative modeling, substituting the default statistical distance of an existing generative adversarial network with CT is shown to consistently improve the performance. PyTorch code is provided.

[Actively Identifying Causal Effects with Latent Variables Given Only Response Variable Observable](#)

- Tian-Zuo Wang, Zhi-Hua Zhou
- abstract@[open-review](#): In many real tasks, it is generally desired to study the causal effect on a specific target (response variable) only, with no need to identify the thorough causal effects involving all variables. In this paper, we attempt to identify such effects by a few active interventions where only the response variable is observable. This task is challenging because the causal graph is unknown and even there may exist latent confounders. To learn the necessary structure for identifying the effects, we provide the graphical characterization that allows us to efficiently estimate all possible causal effects in a partially mixed ancestral graph (PMAG) by generalized back-door criterion. The characterization guides learning a local structure with the interventional data. Theoretical analysis and empirical studies validate the effectiveness and efficiency of our proposed approach.

[Interventional Sum-Product Networks: Causal Inference with Tractable Probabilistic Models](#)

- Matej Zečević, Devendra Dhami, Athresh Karanam, Sriraam Natarajan, Kristian Kersting
- abstract@[open-review](#): While probabilistic models are an important tool for studying causality, doing so suffers from the intractability of inference. As a step towards tractable causal models, we consider the problem of learning interventional distributions using sum-product networks (SPNs) that are over-parameterized by gate functions, e.g., neural networks. Providing an arbitrarily intervened causal graph as input, effectively subsuming Pearl's do-operator, the gate function predicts the parameters of the SPN. The resulting interventional SPNs are motivated and illustrated by a structural causal model themed around personal health. Our empirical evaluation against competing methods from both generative and causal modelling demonstrates that interventional SPNs indeed are both expressive and causally adequate.

[PettingZoo: Gym for Multi-Agent Reinforcement Learning](#)

- J Terry, Benjamin Black, Nathaniel Grammel, Mario Jayakumar, Ananth Hari, Ryan Sullivan, Luis S Santos, Clemens Dieffendahl, Caroline Horsch, Rodrigo Perez-Vicente, Niall Williams, Yashas Lokesh, Praveen Ravi
- abstract@[open-review](#): This paper introduces the PettingZoo library and the accompanying Agent Environment Cycle ("AEC") games model. PettingZoo is a library of diverse sets of multi-agent environments with a universal, elegant Python API. PettingZoo was developed with the goal of accelerating research in Multi-Agent Reinforcement Learning ("MARL"), by making work more interchangeable, accessible and reproducible akin to what OpenAI's Gym library did for single-agent reinforcement learning. PettingZoo's API, while inheriting many features of Gym, is unique amongst MARL APIs in that it's based around the novel AEC games model. We argue, in part through case studies on major problems in popular MARL environments, that the popular game models are poor conceptual models of the games commonly used with MARL, that they promote severe bugs that are hard to detect, and that the AEC games model addresses these problems.

[Parametric Complexity Bounds for Approximating PDEs with Neural Networks](#)

- Tanya Marwah, Zachary Lipton, Andrej Risteski
- abstract@[open-review](#): Recent experiments have shown that deep networks can approximate solutions to high-dimensional PDEs, seemingly escaping the curse of dimensionality. However, questions regarding the theoretical basis for such approximations, including the required network size remain open. In this paper, we investigate the representational power of neural networks for approximating solutions to linear elliptic PDEs with Dirichlet boundary conditions. We prove that when a PDE's coefficients are representable by small neural networks, the parameters required to approximate its solution scale polynomially with the input dimension $\$d\$$ and proportionally to the parameter counts of the coefficient networks. To this end, we develop a proof technique that simulates gradient descent (in an appropriate Hilbert space) by growing a neural network architecture whose iterates each participate as sub-networks in their (slightly larger) successors, and converge to the solution of the PDE.

[Learning-to-learn non-convex piecewise-Lipschitz functions](#)

- Maria-Florina F. Balcan, Mikhail Khodak, Dravyansh Sharma, Ameet Talwalkar
- abstract@[open-review](#): We analyze the meta-learning of the initialization and step-size of learning algorithms for piecewise-Lipschitz functions, a non-convex setting with applications to both machine learning and algorithms. Starting from recent regret bounds for the exponential forecaster on losses with dispersed discontinuities, we generalize them to be initialization-dependent and then use this result to propose a practical meta-learning procedure that learns both the initialization and the step-size of the algorithm from multiple online learning tasks. Asymptotically, we guarantee that the average regret across tasks scales with a natural notion of task-similarity that measures the amount of overlap between near-optimal regions of different tasks. Finally, we instantiate the method and its guarantee in two important settings: robust meta-learning and multi-task data-driven algorithm design.

[Uncertain Decisions Facilitate Better Preference Learning](#)

- Cassidy Laidlaw, Stuart Russell
- abstract@[open-review](#): Existing observational approaches for learning human preferences, such as inverse reinforcement learning, usually make strong assumptions about the observability of the human's environment. However, in reality, people make many important decisions under uncertainty. To better understand preference learning in these cases, we study the setting of inverse decision theory (IDT), a previously proposed framework where a human is observed making non-sequential binary decisions under uncertainty. In IDT, the human's preferences are conveyed through their loss function, which expresses a tradeoff between different types of mistakes. We give the first statistical analysis of IDT, providing conditions necessary to identify these preferences and characterizing the sample complexity—the number of decisions that must be observed to learn the tradeoff the human is making to a desired precision. Interestingly, we show that it is actually easier to identify preferences when the decision problem is more uncertain. Furthermore, uncertain decision problems allow us to relax the unrealistic assumption that the human is an optimal decision maker but still identify their exact preferences; we give sample complexities in this suboptimal case as well. Our analysis contradicts the intuition that partial observability should make preference learning more difficult. It also provides a first step towards understanding and improving preference learning methods for uncertain and suboptimal humans.

[Decision Transformer: Reinforcement Learning via Sequence Modeling](#)

- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, Igor Mordatch
- abstract@[open-review](#): We introduce a framework that abstracts Reinforcement Learning (RL) as a sequence modeling problem. This allows us to draw upon the simplicity and scalability of the Transformer architecture, and associated advances in language modeling such as GPT-x and BERT. In particular, we present Decision Transformer, an architecture that casts the problem of RL as conditional sequence modeling. Unlike prior approaches to RL that fit value functions or compute policy gradients, Decision Transformer simply outputs the optimal actions by leveraging a causally masked Transformer. By conditioning an autoregressive model on the desired return (reward), past states, and actions, our Decision Transformer model can generate future actions that achieve the desired return. Despite its simplicity, Decision Transformer matches or exceeds the performance of state-of-the-art model-free offline RL baselines on Atari, OpenAI Gym, and Key-to-Door tasks.

Probability Paths and the Structure of Predictions over Time

- Zhiyuan Jerry Lin, Hao Sheng, Sharad Goel
- abstract@[open-review](#): In settings ranging from weather forecasts to political prognostications to financial projections, probability estimates of future binary outcomes often evolve over time. For example, the estimated likelihood of rain on a specific day changes by the hour as new information becomes available. Given a collection of such probability paths, we introduce a Bayesian framework -- which we call the Gaussian latent information martingale, or GLIM -- for modeling the structure of dynamic predictions over time. Suppose, for example, that the likelihood of rain in a week is 50%, and consider two hypothetical scenarios. In the first, one expects the forecast to be equally likely to become either 25% or 75% tomorrow; in the second, one expects the forecast to stay constant for the next several days. A time-sensitive decision-maker might select a course of action immediately in the latter scenario, but may postpone their decision in the former, knowing that new information is imminent. We model these trajectories by assuming predictions update according to a latent process of information flow, which is inferred from historical data. In contrast to general methods for time series analysis, this approach preserves important properties of probability paths such as the martingale structure and appropriate amount of volatility and better quantifies future uncertainties around probability paths. We show that GLIM outperforms three popular baseline methods, producing better estimated posterior probability path distributions measured by three different metrics. By elucidating the dynamic structure of predictions over time, we hope to help individuals make more informed choices.

Deep Extended Hazard Models for Survival Analysis

- Qixian Zhong, Jonas W. Mueller, Jane-Ling Wang
- abstract@[open-review](#): Unlike standard prediction tasks, survival analysis requires modeling right censored data, which must be treated with care. While deep neural networks excel in traditional supervised learning, it remains unclear how to best utilize these models in survival analysis. A key question asks which data-generating assumptions of traditional survival models should be retained and which should be made more flexible via the function-approximating capabilities of neural networks. Rather than estimating the survival function targeted by most existing methods, we introduce a Deep Extended Hazard (DeepEH) model to provide a flexible and general framework for deep survival analysis. The extended hazard model includes the conventional Cox proportional hazards and accelerated failure time models as special cases, so DeepEH subsumes the popular Deep Cox proportional hazard (DeepSurv) and Deep Accelerated Failure Time (DeepAFT) models. We additionally provide theoretical support for the proposed DeepEH model by establishing consistency and convergence rate of the survival function estimator, which underscore the attractive feature that deep learning is able to detect low-dimensional structure of data in high-dimensional space. Numerical experiments also provide evidence that the proposed methods outperform existing statistical and deep learning approaches to survival analysis.

TNASP: A Transformer-based NAS Predictor with a Self-evolution Framework

- Shun Lu, Jixiang Li, Jianchao Tan, Sen Yang, Ji Liu
- abstract@[open-review](#): Predictor-based Neural Architecture Search (NAS) continues to be an important topic because it aims to mitigate the time-consuming search procedure of traditional NAS methods. A promising performance predictor determines the quality of final searched models in predictor-based NAS methods. Most existing predictor-based methodologies train model-based predictors under a proxy dataset setting, which may suffer from the accuracy decline and the generalization problem, mainly due to their poor abilities to represent spatial topology information of the graph structure data. Besides the poor encoding for spatial topology information, these works did not take advantage of the temporal information such as historical evaluations during training. Thus, we propose a Transformer-based NAS performance predictor, associated with a Laplacian matrix based positional encoding strategy, which better represents topology information and achieves better performance than previous state-of-the-art methods on NAS-Bench-101, NAS-Bench-201, and DARTS search space. Furthermore, we also propose a self-evolution framework that can fully utilize temporal information as guidance. This framework iteratively involves the evaluations of previously predicted results as constraints into current optimization iteration, thus further improving the performance of our predictor. Such framework is model-agnostic, thus can enhance performance on various backbone structures for the prediction task. Our proposed method helped us rank 2nd among all teams in CVPR 2021 NAS Competition Track 2: Performance Prediction Track.

Automorphic Equivalence-aware Graph Neural Network

- Fengli Xu, Quanming Yao, Pan Hui, Yong Li
- abstract@[open-review](#): Distinguishing the automorphic equivalence of nodes in a graph plays an essential role in many scientific domains, e.g., computational biologist and social network analysis. However, existing graph neural networks (GNNs) fail to capture such an important property. To make GNN aware of automorphic equivalence, we first introduce a localized variant of this concept --- ego-centered automorphic equivalence (Ego-AE). Then, we design a novel variant of GNN, i.e., GRAPE, that uses learnable AE-aware aggregators to explicitly differentiate the Ego-AE of each node's neighbors with the aids of various subgraph templates. While the design of subgraph templates can be hard, we further propose a genetic algorithm to automatically search them from graph data. Moreover, we theoretically prove that GRAPE is expressive in terms of generating distinct representations for nodes with different Ego-AE features, which fills in a fundamental gap of existing GNN variants. Finally, we empirically validate our model on eight real-world graph data, including social network, e-commerce co-purchase network, and citation network, and show that it consistently outperforms existing GNNs. The source code is public available at <https://github.com/tsinghua-fib-lab/GRAPE>.

Random Shuffling Beats SGD Only After Many Epochs on Ill-Conditioned Problems

- Itay Safran, Ohad Shamir
- abstract@[open-review](#): Recently, there has been much interest in studying the convergence rates of without-replacement SGD, and proving that it is faster than with-replacement SGD in the worst case. However, known lower bounds ignore the problem's geometry, including its condition number, whereas the upper bounds explicitly depend on it. Perhaps surprisingly, we prove that when the condition number is taken into account, without-replacement SGD \emph{does not} significantly improve on with-replacement SGD in terms of worst-case bounds, unless the number of epochs (passes over the data) is larger than the condition number. Since many problems in machine learning and other areas are both ill-conditioned and involve large datasets, this indicates that without-replacement does not necessarily improve over with-replacement sampling for realistic iteration budgets. We show this by providing new lower and upper bounds which are tight (up to log factors), for quadratic problems with commuting quadratic terms, precisely quantifying the dependence on the problem parameters.

Analytic Study of Families of Spurious Minima in Two-Layer ReLU Neural Networks: A Tale of Symmetry II

- Yossi Arjevani, Michael Field
- abstract@[open-review](#): We study the optimization problem associated with fitting two-layer ReLU neural networks with respect to the squared loss, where labels are generated by a target network. We make use of the rich symmetry structure to develop a novel set of tools for studying families of spurious minima. In contrast to existing approaches which operate in limiting regimes, our technique directly addresses the nonconvex loss landscape for finite number of inputs d and neurons k , and provides analytic, rather than heuristic, information. In particular, we derive analytic estimates for the loss at different minima, and prove that, modulo $\mathcal{O}(d^{-1/2})$ -terms, the Hessian spectrum concentrates near small positive constants, with the exception of $\Theta(d)$ eigenvalues which grow linearly with d . We further show that the Hessian spectrum at global and spurious minima coincide to $\mathcal{O}(d^{-1/2})$ -order, thus challenging our ability to argue about statistical generalization through local curvature. Lastly, our technique provides the exact \emph{fractional} dimensionality at which families of critical points turn from saddles into spurious minima. This makes possible the study of the creation and the annihilation of spurious minima using powerful tools from equivariant bifurcation theory.

[CAM-GAN: Continual Adaptation Modules for Generative Adversarial Networks](#)

- Sakshi Varshney, Vinay Kumar Verma, P. K. Srijith, Lawrence Carin, Piyush Rai
- abstract@[open-review](#): We present a continual learning approach for generative adversarial networks (GANs), by designing and leveraging parameter-efficient feature map transformations. Our approach is based on learning a set of global and task-specific parameters. The global parameters are fixed across tasks whereas the task-specific parameters act as local adapters for each task, and help in efficiently obtaining task-specific feature maps. Moreover, we propose an element-wise addition of residual bias in the transformed feature space, which further helps stabilize GAN training in such settings. Our approach also leverages task similarities based on the Fisher information matrix. Leveraging this knowledge from previous tasks significantly improves the model performance. In addition, the similarity measure also helps reduce the parameter growth in continual adaptation and helps to learn a compact model. In contrast to the recent approaches for continually-learned GANs, the proposed approach provides a memory-efficient way to perform effective continual data generation. Through extensive experiments on challenging and diverse datasets, we show that the feature-map-transformation approach outperforms state-of-the-art methods for continually-learned GANs, with substantially fewer parameters. The proposed method generates high-quality samples that can also improve the generative-replay-based continual learning for discriminative tasks.

[Structured Dropout Variational Inference for Bayesian Neural Networks](#)

- Son Nguyen, Duong Nguyen, Khai Nguyen, Khoat Than, Hung Bui, Nhat Ho
- abstract@[open-review](#): Approximate inference in Bayesian deep networks exhibits a dilemma of how to yield high fidelity posterior approximations while maintaining computational efficiency and scalability. We tackle this challenge by introducing a novel variational structured approximation inspired by the Bayesian interpretation of Dropout regularization. Concretely, we focus on the inflexibility of the factorized structure in Dropout posterior and then propose an improved method called Variational Structured Dropout (VSD). VSD employs an orthogonal transformation to learn a structured representation on the variational Gaussian noise with plausible complexity, and consequently induces statistical dependencies in the approximate posterior. Theoretically, VSD successfully addresses the pathologies of previous Variational Dropout methods and thus offers a standard Bayesian justification. We further show that VSD induces an adaptive regularization term with several desirable properties which contribute to better generalization. Finally, we conduct extensive experiments on standard benchmarks to demonstrate the effectiveness of VSD over state-of-the-art variational methods on predictive accuracy, uncertainty estimation, and out-of-distribution detection.

[Neural Relightable Participating Media Rendering](#)

- Quan Zheng, Gurprit Singh, Hans-peter Seidel
- abstract@[open-review](#): Learning neural radiance fields of a scene has recently allowed realistic novel view synthesis of the scene, but they are limited to synthesize images under the original fixed lighting condition. Therefore, they are not flexible for the eagerly desired tasks like relighting, scene editing and scene composition. To tackle this problem, several recent methods propose to disentangle reflectance and illumination from the radiance field. These methods can cope with solid objects with opaque surfaces but participating media are neglected. Also, they take into account only direct illumination or at most one-bounce indirect illumination, thus suffer from energy loss due to ignoring the high-order indirect illumination. We propose to learn neural representations for participating media with a complete simulation of global illumination. We estimate direct illumination via ray tracing and compute indirect illumination with spherical harmonics. Our approach avoids computing the lengthy indirect bounces and does not suffer from energy loss. Our experiments on multiple scenes show that our approach achieves superior visual quality and numerical performance compared to state-of-the-art methods, and it can generalize to deal with solid objects with opaque surfaces as well.

[Efficient Neural Network Training via Forward and Backward Propagation Sparsification](#)

- Xiao Zhou, Weizhong Zhang, Zonghao Chen, SHIZHE DIAO, Tong Zhang
- abstract@[open-review](#): Sparse training is a natural idea to accelerate the training speed of deep neural networks and save the memory usage, especially since large modern neural networks are significantly over-parameterized. However, most of the existing methods cannot achieve this goal in practice because the chain rule based gradient (w.r.t. structure parameters) estimators adopted by previous methods require dense computation at least in the backward propagation step. This paper solves this problem by proposing an efficient sparse training method with completely sparse forward and backward passes. We first formulate the training process as a continuous minimization problem under global sparsity constraint. We then separate the optimization process into two steps, corresponding to weight update and structure parameter update. For the former step, we use the conventional chain rule, which can be sparse via exploiting the sparse structure. For the latter step, instead of using the chain rule based gradient estimators as in existing methods, we propose a variance reduced policy gradient estimator, which only requires two forward passes without backward propagation, thus achieving completely sparse training. We prove that the variance of our gradient estimator is bounded. Extensive experimental results on real-world datasets demonstrate that compared to previous methods, our algorithm is much more effective in accelerating the training process, up to an order of magnitude faster.

[Learning to Ground Multi-Agent Communication with Autoencoders](#)

- Toru Lin, Jacob Huh, Christopher Stauffer, Ser Nam Lim, Phillip Isola
- abstract@[open-review](#): Communication requires having a common language, a lingua franca, between agents. This language could emerge via a consensus process, but it may require many generations of trial and error. Alternatively, the lingua franca can be given by the environment, where agents ground their language in representations of the observed world. We demonstrate a simple way to ground language in learned representations, which facilitates decentralized multi-agent communication and coordination. We find that a standard representation learning algorithm -- autoencoding -- is sufficient for arriving at a grounded common language. When agents broadcast these representations, they learn to understand and respond to each other's utterances and achieve surprisingly strong task performance across a variety of multi-agent communication environments.

[Large-Scale Wasserstein Gradient Flows](#)

- Petr Mokrov, Alexander Korotin, Lingxiao Li, Aude Genevay, Justin M. Solomon, Evgeny Burnaev
- abstract@[open-review](#): Wasserstein gradient flows provide a powerful means of understanding and solving many diffusion equations. Specifically, Fokker-Planck equations, which model the diffusion of probability measures, can be understood as gradient descent over entropy functionals in Wasserstein space. This equivalence, introduced by Jordan, Kinderlehrer and Otto, inspired the so-called JKO scheme to approximate these diffusion processes via an implicit discretization of the gradient flow in Wasserstein space. Solving the optimization problem associated with each JKO step, however, presents serious computational challenges. We introduce a scalable method to approximate Wasserstein gradient flows, targeted to machine learning applications. Our approach relies on input-convex neural networks (ICNNs) to discretize the JKO steps, which can be optimized by stochastic gradient descent. Contrarily to previous work, our method does not require domain discretization or particle simulation. As a result, we can sample from the measure at each time step of the diffusion and compute its probability density. We demonstrate the performance of our algorithm by computing diffusions following the Fokker-Planck equation and apply it to unnormalized density sampling as well as nonlinear filtering.

[Who Leads and Who Follows in Strategic Classification?](#)

- Tijana Zrnic, Eric Mazumdar, Shankar Sastry, Michael Jordan

- abstract@[open-review](#): As predictive models are deployed into the real world, they must increasingly contend with strategic behavior. A growing body of work on strategic classification treats this problem as a Stackelberg game: the decision-maker "leads" in the game by deploying a model, and the strategic agents "follow" by playing their best response to the deployed model. Importantly, in this framing, the burden of learning is placed solely on the decision-maker, while the agents' best responses are implicitly treated as instantaneous. In this work, we argue that the order of play in strategic classification is fundamentally determined by the relative frequencies at which the decision-maker and the agents adapt to each other's actions. In particular, by generalizing the standard model to allow both players to learn over time, we show that a decision-maker that makes updates faster than the agents can reverse the order of play, meaning that the agents lead and the decision-maker follows. We observe in standard learning settings that such a role reversal can be desirable for both the decision-maker and the strategic agents. Finally, we show that a decision-maker with the freedom to choose their update frequency can induce learning dynamics that converge to Stackelberg equilibria with either order of play.

Unadversarial Examples: Designing Objects for Robust Vision

- Hadi Salman, Andrew Ilyas, Logan Engstrom, Sai Venkrala, Aleksander Madry, Ashish Kapoor
- abstract@[open-review](#): We study a class of computer vision settings wherein one can modify the design of the objects being recognized. We develop a framework that leverages this capability---and deep networks' unusual sensitivity to input perturbations---to design ``robust objects," i.e., objects that are explicitly optimized to be confidently classified. Our framework yields improved performance on standard benchmarks, a simulated robotics environment, and physical-world experiments.

Deep Jump Learning for Off-Policy Evaluation in Continuous Treatment Settings

- Hengrui Cai, Chengchun Shi, Rui Song, Wenbin Lu
- abstract@[open-review](#): We consider off-policy evaluation (OPE) in continuous treatment settings, such as personalized dose-finding. In OPE, one aims to estimate the mean outcome under a new treatment decision rule using historical data generated by a different decision rule. Most existing works on OPE focus on discrete treatment settings. To handle continuous treatments, we develop a novel estimation method for OPE using deep jump learning. The key ingredient of our method lies in adaptively discretizing the treatment space using deep discretization, by leveraging deep learning and multi-scale change point detection. This allows us to apply existing OPE methods in discrete treatments to handle continuous treatments. Our method is further justified by theoretical results, simulations, and a real application to Warfarin Dosing.

Attention Approximates Sparse Distributed Memory

- Trenton Bricken, Cengiz Pehlevan
- abstract@[open-review](#): While Attention has come to be an important mechanism in deep learning, there remains limited intuition for why it works so well. Here, we show that Transformer Attention can be closely related under certain data conditions to Kanerva's Sparse Distributed Memory (SDM), a biologically plausible associative memory model. We confirm that these conditions are satisfied in pre-trained GPT2 Transformer models. We discuss the implications of the Attention-SDM map and provide new computational and biological interpretations of Attention.

Augmented Shortcuts for Vision Transformers

- Yehui Tang, Kai Han, Chang Xu, An Xiao, Yiping Deng, Chao Xu, Yunhe Wang
- abstract@[open-review](#): Transformer models have achieved great progress on computer vision tasks recently. The rapid development of vision transformers is mainly contributed by their high representation ability for extracting informative features from input images. However, the mainstream transformer models are designed with deep architectures, and the feature diversity will be continuously reduced as the depth increases, \ie, feature collapse. In this paper, we theoretically analyze the feature collapse phenomenon and study the relationship between shortcuts and feature diversity in these transformer models. Then, we present an augmented shortcut scheme, which inserts additional paths with learnable parameters in parallel on the original shortcuts. To save the computational costs, we further explore an efficient approach that uses the block-circulant projection to implement augmented shortcuts. Extensive experiments conducted on benchmark datasets demonstrate the effectiveness of the proposed method, which brings about 1% accuracy increase of the state-of-the-art visual transformers without obviously increasing their parameters and FLOPs.

Finding Regions of Heterogeneity in Decision-Making via Expected Conditional Covariance

- Justin Lim, Christina X Ji, Michael Oberst, Saul Blecker, Leora Horwitz, David Sontag
- abstract@[open-review](#): Individuals often make different decisions when faced with the same context, due to personal preferences and background. For instance, judges may vary in their leniency towards certain drug-related offenses, and doctors may vary in their preference for how to start treatment for certain types of patients. With these examples in mind, we present an algorithm for identifying types of contexts (e.g., types of cases or patients) with high inter-decision-maker disagreement. We formalize this as a causal inference problem, seeking a region where the assignment of decision-maker has a large causal effect on the decision. Our algorithm finds such a region by maximizing an empirical objective, and we give a generalization bound for its performance. In a semi-synthetic experiment, we show that our algorithm recovers the correct region of heterogeneity accurately compared to baselines. Finally, we apply our algorithm to real-world healthcare datasets, recovering variation that aligns with existing clinical knowledge.

Identifying and Benchmarking Natural Out-of-Context Prediction Problems

- David Madras, Richard Zemel
- abstract@[open-review](#): Deep learning systems frequently fail at out-of-context (OOC) prediction, the problem of making reliable predictions on uncommon or unusual inputs or subgroups of the training distribution. To this end, a number of benchmarks for measuring OOC performance have been recently introduced. In this work, we introduce a framework unifying the literature on OOC performance measurement, and demonstrate how rich auxiliary information can be leveraged to identify candidate sets of OOC examples in existing datasets. We present NOOCh: a suite of naturally-occurring "challenge sets", and show how varying notions of context can be used to probe specific OOC failure modes. Experimentally, we explore the tradeoffs between various learning approaches on these challenge sets and demonstrate how the choices made in designing OOC benchmarks can yield varying conclusions.

Label Disentanglement in Partition-based Extreme Multilabel Classification

- Xuanqing Liu, Wei-Cheng Chang, Hsiang-Fu Yu, Cho-Jui Hsieh, Inderjit Dhillon
- abstract@[open-review](#): Partition-based methods are increasingly-used in extreme multi-label classification (XMC) problems due to their scalability to large output spaces (e.g., millions or more). However, existing methods partition the large label space into mutually exclusive clusters, which is sub-optimal when labels have multi-modality and rich semantics. For instance, the label "Apple" can be the fruit or the brand name, which leads to the following research question: can we disentangle these multi-modal labels with non-exclusive clustering tailored for downstream XMC tasks? In this paper, we show that the label assignment problem in partition-based XMC can be formulated as an optimization problem, with the objective of maximizing precision rates. This leads to an efficient algorithm to form flexible and overlapped label clusters, and a method that can alternatively optimize the cluster assignments and the model parameters for partition-based XMC. Experimental results on synthetic and real datasets show that our method can successfully disentangle multi-modal labels, leading to state-of-the-art (SOTA) results on four XMC benchmarks.

[Leveraging SE\(3\) Equivariance for Self-supervised Category-Level Object Pose Estimation from Point Clouds](#)

- Xiaolong Li, Yijia Weng, Li Yi, Leonidas J. Guibas, A. Abbott, Shuran Song, He Wang
- abstract@[open-review](#): Category-level object pose estimation aims to find 6D object poses of previously unseen object instances from known categories without access to object CAD models. To reduce the huge amount of pose annotations needed for category-level learning, we propose for the first time a self-supervised learning framework to estimate category-level 6D object pose from single 3D point clouds. During training, our method assumes no ground-truth pose annotations, no CAD models, and no multi-view supervision. The key to our method is to disentangle shape and pose through an invariant shape reconstruction module and an equivariant pose estimation module, empowered by SE(3) equivariant point cloud networks. The invariant shape reconstruction module learns to perform aligned reconstructions, yielding a category-level reference frame without using any annotations. In addition, the equivariant pose estimation module achieves category-level pose estimation accuracy that is comparable to some fully supervised methods. Extensive experiments demonstrate the effectiveness of our approach on both complete and partial depth point clouds from the ModelNet40 benchmark, and on real depth point clouds from the NOCS-REAL 275 dataset. The project page with code and visualizations can be found at: dragonlong.github.io/equi-pose.

[A Theoretical Analysis of Fine-tuning with Linear Teachers](#)

- Gal Shachaf, Alon Brutzkus, Amir Globerson
- abstract@[open-review](#): Fine-tuning is a common practice in deep learning, achieving excellent generalization results on downstream tasks using relatively little training data. Although widely used in practice, it is not well understood theoretically. Here we analyze the sample complexity of this scheme for regression with linear teachers in several settings. Intuitively, the success of fine-tuning depends on the similarity between the source tasks and the target task. But what is the right way of measuring this similarity? We show that the relevant measure has to do with the relation between the source task, the target task and the covariance structure of the target data. In the setting of linear regression, we show that under realistic settings there can be substantial sample complexity reduction when the above measure is low. For deep linear regression, we propose a novel result regarding the inductive bias of gradient-based training when the network is initialized with pretrained weights. Using this result we show that the similarity measure for this setting is also affected by the depth of the network. We conclude with results on shallow ReLU models, and analyze the dependence of sample complexity there on source and target tasks. We empirically demonstrate our results for both synthetic and realistic data.

[Overinterpretation reveals image classification model pathologies](#)

- Brandon Carter, Siddhartha Jain, Jonas W. Mueller, David Gifford
- abstract@[open-review](#): Image classifiers are typically scored on their test set accuracy, but high accuracy can mask a subtle type of model failure. We find that high scoring convolutional neural networks (CNNs) on popular benchmarks exhibit troubling pathologies that allow them to display high accuracy even in the absence of semantically salient features. When a model provides a high-confidence decision without salient supporting input features, we say the classifier has overinterpreted its input, finding too much class-evidence in patterns that appear nonsensical to humans. Here, we demonstrate that neural networks trained on CIFAR-10 and ImageNet suffer from overinterpretation, and we find models on CIFAR-10 make confident predictions even when 95% of input images are masked and humans cannot discern salient features in the remaining pixel-subsets. We introduce Batched Gradient SIS, a new method for discovering sufficient input subsets for complex datasets, and use this method to show the sufficiency of border pixels in ImageNet for training and testing. Although these patterns portend potential model fragility in real-world deployment, they are in fact valid statistical patterns of the benchmark that alone suffice to attain high test accuracy. Unlike adversarial examples, overinterpretation relies upon unmodified image pixels. We find ensembling and input dropout can each help mitigate overinterpretation.

[Neural Circuit Synthesis from Specification Patterns](#)

- Frederik Schmitt, Christopher Hahn, Markus N Rabe, Bernd Finkbeiner
- abstract@[open-review](#): We train hierarchical Transformers on the task of synthesizing hardware circuits directly out of high-level logical specifications in linear-time temporal logic (LTL). The LTL synthesis problem is a well-known algorithmic challenge with a long history and an annual competition is organized to track the improvement of algorithms and tooling over time. New approaches using machine learning might open a lot of possibilities in this area, but suffer from the lack of sufficient amounts of training data. In this paper, we consider a method to generate large amounts of additional training data, i.e., pairs of specifications and circuits implementing them. We ensure that this synthetic data is sufficiently close to human-written specifications by mining common patterns from the specifications used in the synthesis competitions. We show that hierarchical Transformers trained on this synthetic data solve a significant portion of problems from the synthesis competitions, and even out-of-distribution examples from a recent case study.

[Directional Message Passing on Molecular Graphs via Synthetic Coordinates](#)

- Johannes Gasteiger, Chandan Yeshwanth, Stephan Günnemann
- abstract@[open-review](#): Graph neural networks that leverage coordinates via directional message passing have recently set the state of the art on multiple molecular property prediction tasks. However, they rely on atom position information that is often unavailable, and obtaining it is usually prohibitively expensive or even impossible. In this paper we propose synthetic coordinates that enable the use of advanced GNNs without requiring the true molecular configuration. We propose two distances as synthetic coordinates: Distance bounds that specify the rough range of molecular configurations, and graph-based distances using a symmetric variant of personalized PageRank. To leverage both distance and angular information we propose a method of transforming normal graph neural networks into directional MPNNs. We show that with this transformation we can reduce the error of a normal graph neural network by 55% on the ZINC benchmark. We furthermore set the state of the art on ZINC and coordinate-free QM9 by incorporating synthetic coordinates in the SMP and DimeNet++ models. Our implementation is available online.

[Federated Multi-Task Learning under a Mixture of Distributions](#)

- Othmane Marfoq, Giovanni Neglia, Aurélien Bellet, Laetitia Kameni, Richard Vidal
- abstract@[open-review](#): The increasing size of data generated by smartphones and IoT devices motivated the development of Federated Learning (FL), a framework for on-device collaborative training of machine learning models. First efforts in FL focused on learning a single global model with good average performance across clients, but the global model may be arbitrarily bad for a given client, due to the inherent heterogeneity of local data distributions. Federated multi-task learning (MTL) approaches can learn personalized models by formulating an opportune penalized optimization problem. The penalization term can capture complex relations among personalized models, but eschews clear statistical assumptions about local data distributions. In this work, we propose to study federated MTL under the flexible assumption that each local data distribution is a mixture of unknown underlying distributions. This assumption encompasses most of the existing personalized FL approaches and leads to federated EM-like algorithms for both client-server and fully decentralized settings. Moreover, it provides a principled way to serve personalized models to clients not seen at training time. The algorithms' convergence is analyzed through a novel federated surrogate optimization framework, which can be of general interest. Experimental results on FL benchmarks show that our approach provides models with higher accuracy and fairness than state-of-the-art methods.

[Learning Generative Vision Transformer with Energy-Based Latent Space for Saliency Prediction](#)

- Jing Zhang, Jianwen Xie, Nick Barnes, Ping Li

- abstract@[open-review](#): Vision transformer networks have shown superiority in many computer vision tasks. In this paper, we take a step further by proposing a novel generative vision transformer with latent variables following an informative energy-based prior for salient object detection. Both the vision transformer network and the energy-based prior model are jointly trained via Markov chain Monte Carlo-based maximum likelihood estimation, in which the sampling from the intractable posterior and prior distributions of the latent variables are performed by Langevin dynamics. Further, with the generative vision transformer, we can easily obtain a pixel-wise uncertainty map from an image, which indicates the model confidence in predicting saliency from the image. Different from the existing generative models which define the prior distribution of the latent variables as a simple isotropic Gaussian distribution, our model uses an energy-based informative prior which can be more expressive to capture the latent space of the data. We apply the proposed framework to both RGB and RGB-D salient object detection tasks. Extensive experimental results show that our framework can achieve not only accurate saliency predictions but also meaningful uncertainty maps that are consistent with the human perception.

[Regularization in ResNet with Stochastic Depth](#)

- Soufiane Hayou, Fadhel Ayed
- abstract@[open-review](#): Regularization plays a major role in modern deep learning. From classic techniques such as L1, L2 penalties to other noise-based methods such as Dropout, regularization often yields better generalization properties by avoiding overfitting. Recently, Stochastic Depth (SD) has emerged as an alternative regularization technique for residual neural networks (ResNets) and has proven to boost the performance of ResNet on many tasks [Huang et al., 2016]. Despite the recent success of SD, little is known about this technique from a theoretical perspective. This paper provides a hybrid analysis combining perturbation analysis and signal propagation to shed light on different regularization effects of SD. Our analysis allows us to derive principled guidelines for choosing the survival rates used for training with SD.

[ResT: An Efficient Transformer for Visual Recognition](#)

- Qinglong Zhang, Yu-Bin Yang
- abstract@[open-review](#): This paper presents an efficient multi-scale vision Transformer, called ResT, that capably served as a general-purpose backbone for image recognition. Unlike existing Transformer methods, which employ standard Transformer blocks to tackle raw images with a fixed resolution, our ResT have several advantages: (1) A memory-efficient multi-head self-attention is built, which compresses the memory by a simple depth-wise convolution, and projects the interaction across the attention-heads dimension while keeping the diversity ability of multi-heads; (2) Positional encoding is constructed as spatial attention, which is more flexible and can tackle with input images of arbitrary size without interpolation or fine-tune; (3) Instead of the straightforward tokenization at the beginning of each stage, we design the patch embedding as a stack of overlapping convolution operation with stride on the token map. We comprehensively validate ResT on image classification and downstream tasks. Experimental results show that the proposed ResT can outperform the recently state-of-the-art backbones by a large margin, demonstrating the potential of ResT as strong backbones. The code and models will be made publicly available at <https://github.com/wofmanaf/ResT>.

[Adversarial Examples for k-Nearest Neighbor Classifiers Based on Higher-Order Voronoi Diagrams](#)

- Chawin Sitawarin, Evgenios Kornaropoulos, Dawn Song, David Wagner
- abstract@[open-review](#): Adversarial examples are a widely studied phenomenon in machine learning models. While most of the attention has been focused on neural networks, other practical models also suffer from this issue. In this work, we propose an algorithm for evaluating the adversarial robustness of $\$k\$$ -nearest neighbor classification, i.e., finding a minimum-norm adversarial example. Diverging from previous proposals, we propose the first geometric approach by performing a search that expands outwards from a given input point. On a high level, the search radius expands to the nearby higher-order Voronoi cells until we find a cell that classifies differently from the input point. To scale the algorithm to a large $\$k$$, we introduce approximation steps that find perturbation with smaller norm, compared to the baselines, in a variety of datasets. Furthermore, we analyze the structural properties of a dataset where our approach outperforms the competition.

[Adversarially Robust 3D Point Cloud Recognition Using Self-Supervisions](#)

- Jiachen Sun, Yulong Cao, Christopher B Choy, Zhiding Yu, Anima Anandkumar, Zhuoqing Morley Mao, Chaowei Xiao
- abstract@[open-review](#): 3D point cloud data is increasingly used in safety-critical applications such as autonomous driving. Thus, the robustness of 3D deep learning models against adversarial attacks becomes a major consideration. In this paper, we systematically study the impact of various self-supervised learning proxy tasks on different architectures and threat models for 3D point clouds with adversarial training. Specifically, we study MLP-based (PointNet), convolution-based (DGCNN), and transformer-based (PCT) 3D architectures. Through extensive experimentation, we demonstrate that appropriate applications of self-supervision can significantly enhance the robustness in 3D point cloud recognition, achieving considerable improvements compared to the standard adversarial training baseline. Our analysis reveals that local feature learning is desirable for adversarial robustness in point clouds since it limits the adversarial propagation between the point-level input perturbations and the model's final output. This insight also explains the success of DGCNN and the jigsaw proxy task in achieving stronger 3D adversarial robustness.

[Tuning Mixed Input Hyperparameters on the Fly for Efficient Population Based AutoRL](#)

- Jack Parker-Holder, Vu Nguyen, Shaan Desai, Stephen J Roberts
- abstract@[open-review](#): Despite a series of recent successes in reinforcement learning (RL), many RL algorithms remain sensitive to hyperparameters. As such, there has recently been interest in the field of AutoRL, which seeks to automate design decisions to create more general algorithms. Recent work suggests that population based approaches may be effective AutoRL algorithms, by learning hyperparameter schedules on the fly. In particular, the PB2 algorithm is able to achieve strong performance in RL tasks by formulating online hyperparameter optimization as time varying GP-bandit problem, while also providing theoretical guarantees. However, PB2 is only designed to work for \emph{continuous} hyperparameters, which severely limits its utility in practice. In this paper we introduce a new (provably) efficient hierarchical approach for optimizing \emph{both continuous and categorical} variables, using a new time-varying bandit algorithm specifically designed for the population based training regime. We evaluate our approach on the challenging Procgen benchmark, where we show that explicitly modelling dependence between data augmentation and other hyperparameters improves generalization.

[Neural Algorithmic Reasoners are Implicit Planners](#)

- Andreea-Ioana Deac, Petar Veličković, Ognjen Milinkovic, Pierre-Luc Bacon, Jian Tang, Mladen Nikolic
- abstract@[open-review](#): Implicit planning has emerged as an elegant technique for combining learned models of the world with end-to-end model-free reinforcement learning. We study the class of implicit planners inspired by value iteration, an algorithm that is guaranteed to yield perfect policies in fully-specified tabular environments. We find that prior approaches either assume that the environment is provided in such a tabular form---which is highly restrictive---or infer "local neighbourhoods" of states to run value iteration over---for which we discover an algorithmic bottleneck effect. This effect is caused by explicitly running the planning algorithm based on scalar predictions in every state, which can be harmful to data efficiency if such scalars are improperly predicted. We propose eXecuted Latent Value Iteration Networks (XLVINs), which alleviate the above limitations. Our method performs all planning computations in a high-dimensional latent space, breaking the algorithmic bottleneck. It maintains alignment with value iteration by carefully leveraging neural graph-algorithmic reasoning and contrastive self-supervised learning. Across seven low-data settings---including classical control, navigation and Atari---XLVINs provide significant improvements to data efficiency against value iteration-based implicit planners, as well as relevant model-free baselines. Lastly, we empirically verify that XLVINs can closely align with value iteration.

Self-Supervised Learning with Kernel Dependence Maximization

- Yazhe Li, Roman Pogodin, Danica J. Sutherland, Arthur Gretton
- abstract@[open-review](#): We approach self-supervised learning of image representations from a statistical dependence perspective, proposing Self-Supervised Learning with the Hilbert-Schmidt Independence Criterion (SSL-HSIC). SSL-HSIC maximizes dependence between representations of transformations of an image and the image identity, while minimizing the kernelized variance of those representations. This framework yields a new understanding of InfoNCE, a variational lower bound on the mutual information (MI) between different transformations. While the MI itself is known to have pathologies which can result in learning meaningless representations, its bound is much better behaved: we show that it implicitly approximates SSL-HSIC (with a slightly different regularizer). Our approach also gives us insight into BYOL, a negative-free SSL method, since SSL-HSIC similarly learns local neighborhoods of samples. SSL-HSIC allows us to directly optimize statistical dependence in time linear in the batch size, without restrictive data assumptions or indirect mutual information estimators. Trained with or without a target network, SSL-HSIC matches the current state-of-the-art for standard linear evaluation on ImageNet, semi-supervised learning and transfer to other classification and vision tasks such as semantic segmentation, depth estimation and object recognition. Code is available at https://github.com/deepmind/ssl_hsic.

CROCS: Clustering and Retrieval of Cardiac Signals Based on Patient Disease Class, Sex, and Age

- Dani Kiyasseh, Tingting Zhu, David Clifton
- abstract@[open-review](#): The process of manually searching for relevant instances in, and extracting information from, clinical databases underpin a multitude of clinical tasks. Such tasks include disease diagnosis, clinical trial recruitment, and continuing medical education. This manual search-and-extract process, however, has been hampered by the growth of large-scale clinical databases and the increased prevalence of unlabelled instances. To address this challenge, we propose a supervised contrastive learning framework, CROCS, where representations of cardiac signals associated with a set of patient-specific attributes (e.g., disease class, sex, age) are attracted to learnable embeddings entitled clinical prototypes. We exploit such prototypes for both the clustering and retrieval of unlabelled cardiac signals based on multiple patient attributes. We show that CROCS outperforms the state-of-the-art method, DTC, when clustering and also retrieves relevant cardiac signals from a large database. We also show that clinical prototypes adopt a semantically meaningful arrangement based on patient attributes and thus confer a high degree of interpretability.

Representing Hyperbolic Space Accurately using Multi-Component Floats

- Tao Yu, Christopher M. De Sa
- abstract@[open-review](#): Hyperbolic space is particularly useful for embedding data with hierarchical structure; however, representing hyperbolic space with ordinary floating-point numbers greatly affects the performance due to its \emph{ineluctable} numerical errors. Simply increasing the precision of floats fails to solve the problem and incurs a high computation cost for simulating greater-than-double-precision floats on hardware such as GPUs, which does not support them. In this paper, we propose a simple, feasible-on-GPUs, and easy-to-understand solution for numerically accurate learning on hyperbolic space. We do this with a new approach to represent hyperbolic space using multi-component floating-point (MCF) in the Poincar\'e upper-half space model. Theoretically and experimentally we show our model has small numerical error, and on embedding tasks across various datasets, models represented by multi-component floating-points gain more capacity and run significantly faster on GPUs than prior work.

Dimensionality Reduction for Wasserstein Barycenter

- Zachary Izzo, Sandeep Silwal, Samson Zhou
- abstract@[open-review](#): The Wasserstein barycenter is a geometric construct which captures the notion of centrality among probability distributions, and which has found many applications in machine learning. However, most algorithms for finding even an approximate barycenter suffer an exponential dependence on the dimension d of the underlying space of the distributions. In order to cope with this ``curse of dimensionality," we study dimensionality reduction techniques for the Wasserstein barycenter problem. When the barycenter is restricted to support of size n , we show that randomized dimensionality reduction can be used to map the problem to a space of dimension $O(\log n)$ independent of both d and k , and that \emph{any} solution found in the reduced dimension will have its cost preserved up to arbitrary small error in the original space. We provide matching upper and lower bounds on the size of the reduced dimension, showing that our methods are optimal up to constant factors. We also provide a coresnet construction for the Wasserstein barycenter problem that significantly decreases the number of input distributions. The coresets can be used in conjunction with random projections and thus further improve computation time. Lastly, our experimental results validate the speedup provided by dimensionality reduction while maintaining solution quality.

Neural Population Geometry Reveals the Role of Stochasticity in Robust Perception

- Joel Dapello, Jenelle Feather, Hang Le, Tiago Marques, David Cox, Josh McDermott, James J DiCarlo, Sueyeon Chung
- abstract@[open-review](#): Adversarial examples are often cited by neuroscientists and machine learning researchers as an example of how computational models diverge from biological sensory systems. Recent work has proposed adding biologically-inspired components to visual neural networks as a way to improve their adversarial robustness. One surprisingly effective component for reducing adversarial vulnerability is response stochasticity, like that exhibited by biological neurons. Here, using recently developed geometrical techniques from computational neuroscience, we investigate how adversarial perturbations influence the internal representations of standard, adversarially trained, and biologically-inspired stochastic networks. We find distinct geometric signatures for each type of network, revealing different mechanisms for achieving robust representations. Next, we generalize these results to the auditory domain, showing that neural stochasticity also makes auditory models more robust to adversarial perturbations. Geometric analysis of the stochastic networks reveals overlap between representations of clean and adversarially perturbed stimuli, and quantitatively demonstrate that competing geometric effects of stochasticity mediate a tradeoff between adversarial and clean performance. Our results shed light on the strategies of robust perception utilized by adversarially trained and stochastic networks, and help explain how stochasticity may be beneficial to machine and biological computation.

Unsupervised Learning of Compositional Energy Concepts

- Yilun Du, Shuang Li, Yash Sharma, Josh Tenenbaum, Igor Mordatch
- abstract@[open-review](#): Humans are able to rapidly understand scenes by utilizing concepts extracted from prior experience. Such concepts are diverse, and include global scene descriptors, such as the weather or lighting, as well as local scene descriptors, such as the color or size of a particular object. So far, unsupervised discovery of concepts has focused on either modeling the global scene-level or the local object-level factors of variation, but not both. In this work, we propose COMET, which discovers and represents concepts as separate energy functions, enabling us to represent both global concepts as well as objects under a unified framework. COMET discovers energy functions through recomposing the input image, which we find captures independent factors without additional supervision. Sample generation in COMET is formulated as an optimization process on underlying energy functions, enabling us to generate images with permuted and composed concepts. Finally, discovered visual concepts in COMET generalize well, enabling us to compose concepts between separate modalities of images as well as with other concepts discovered by a separate instance of COMET trained on a different dataset. Code and data available at <https://energy-based-model.github.io/comet/>.

Nearly Horizon-Free Offline Reinforcement Learning

- Tongzheng Ren, Jialian Li, Bo Dai, Simon S. Du, Sujay Sanghavi
- abstract@[open-review](#): We revisit offline reinforcement learning on episodic time-homogeneous Markov Decision Processes (MDP). For tabular MDP with $\$S\$$ states and $\$A\$$ actions, or linear MDP with anchor points and feature dimension $\$d\$$, given the collected $\$K\$$ episodes data with minimum visiting probability of (anchor) state-action pairs $\$d_m\$$, we obtain nearly horizon $\$H\$$ -free sample complexity bounds for offline reinforcement learning when the total reward is upper bounded by 1. Specifically:
 - For offline policy evaluation, we obtain an $\$tilde{O}(\sqrt{\frac{1}{Kd_m}})\$$ error bound for the plug-in estimator, which matches the lower bound up to logarithmic factors and does not have additional dependency on $\$mathrm{poly}(H, S, A, d)\$$ in higher-order term.
 - For offline policy optimization, we obtain an $\$tilde{O}(\sqrt{\frac{1}{Kd_m}} + \frac{\min(S, d)}{Kd_m})\$$ sub-optimality gap for the empirical optimal policy, which approaches the lower bound up to logarithmic factors and a high-order term, improving upon the best known result by [Cui and Yang 2020] that has additional $\$mathrm{poly}(H, S, d)\$$ factors in the main term.To the best of our knowledge, these are the first set of nearly horizon-free bounds for episodic time-homogeneous offline tabular MDP and linear MDP with anchor points. Central to our analysis is a simple yet effective recursion based method to bound a "total variance" term in the offline scenarios, which could be of individual interest.

[Combinatorial Optimization for Panoptic Segmentation: A Fully Differentiable Approach](#)

- Ahmed Abbas, Paul Swoboda
- abstract@[open-review](#): We propose a fully differentiable architecture for simultaneous semantic and instance segmentation (a.k.a. panoptic segmentation) consisting of a convolutional neural network and an asymmetric multiway cut problem solver. The latter solves a combinatorial optimization problem that elegantly incorporates semantic and boundary predictions to produce a panoptic labeling. Our formulation allows to directly maximize a smooth surrogate of the panoptic quality metric by backpropagating the gradient through the optimization problem. Experimental evaluation shows improvement by backpropagating through the optimization problem w.r.t. comparable approaches on Cityscapes and COCO datasets. Overall, our approach of combinatorial optimization for panoptic segmentation (COPS) shows the utility of using optimization in tandem with deep learning in a challenging large scale real-world problem and showcases benefits and insights into training such an architecture.

[Reinforcement Learning with State Observation Costs in Action-Contingent Noiselessly Observable Markov Decision Processes](#)

- HyunJi Alex Nam, Scott Fleming, Emma Brunskill
- abstract@[open-review](#): Many real-world problems that require making optimal sequences of decisions under uncertainty involve costs when the agent wishes to obtain information about its environment. We design and analyze algorithms for reinforcement learning (RL) in Action-Contingent Noiselessly Observable MDPs (ACNO-MDPs), a special class of POMDPs in which the agent can choose to either (1) fully observe the state at a cost and then act; or (2) act without any immediate observation information, relying on past observations to infer the underlying state. ACNO-MDPs arise frequently in important real-world application domains like healthcare, in which clinicians must balance the value of information gleaned from medical tests (e.g., blood-based biomarkers) with the costs of gathering that information (e.g., the costs of labor and materials required to administer such tests). We develop a PAC RL algorithm for tabular ACNO-MDPs that provides substantially tighter bounds, compared to generic POMDP-RL algorithms, on the total number of episodes exhibiting worse than near-optimal performance. For continuous-state ACNO-MDPs, we propose a novel method of incorporating observation information that, when coupled with modern RL algorithms, yields significantly faster learning compared to other POMDP-RL algorithms in several simulated environments.

[Iterative Amortized Policy Optimization](#)

- Joseph Marino, Alexandre Piche, Alessandro Davide Ialongo, Yisong Yue
- abstract@[open-review](#): Policy networks are a central feature of deep reinforcement learning (RL) algorithms for continuous control, enabling the estimation and sampling of high-value actions. From the variational inference perspective on RL, policy networks, when used with entropy or KL regularization, are a form of amortized optimization, optimizing network parameters rather than the policy distributions directly. However, direct amortized mappings can yield suboptimal policy estimates and restricted distributions, limiting performance and exploration. Given this perspective, we consider the more flexible class of iterative amortized optimizers. We demonstrate that the resulting technique, iterative amortized policy optimization, yields performance improvements over direct amortization on benchmark continuous control tasks.

[Revisiting the Calibration of Modern Neural Networks](#)

- Matthias Minderer, Josip Djolonga, Rob Romijnders, Frances Hubis, Xiaohua Zhai, Neil Houlsby, Dustin Tran, Mario Lucic
- abstract@[open-review](#): Accurate estimation of predictive uncertainty (model calibration) is essential for the safe application of neural networks. Many instances of miscalibration in modern neural networks have been reported, suggesting a trend that newer, more accurate models produce poorly calibrated predictions. Here, we revisit this question for recent state-of-the-art image classification models. We systematically relate model calibration and accuracy, and find that the most recent models, notably those not using convolutions, are among the best calibrated. Trends observed in prior model generations, such as decay of calibration with distribution shift or model size, are less pronounced in recent architectures. We also show that model size and amount of pretraining do not fully explain these differences, suggesting that architecture is a major determinant of calibration properties.

[The decomposition of the higher-order homology embedding constructed from the \$\\$k\\$\$ -Laplacian](#)

- Yu-Chia Chen, Marina Meila
- abstract@[open-review](#): The null space of the $\$k\$$ -th order Laplacian $\$mathbf{\mathcal{L}}_k\$$, known as the $\{\text{em } \$k\$-th homology vector space\}$, encodes the non-trivial topology of a manifold or a network. Understanding the structure of the homology embedding can thus disclose geometric or topological information from the data. The study of the null space embedding of the graph Laplacian $\$mathbf{\mathcal{L}}_0\$$ has spurred new research and applications, such as spectral clustering algorithms with theoretical guarantees and estimators of the Stochastic Block Model. In this work, we investigate the geometry of the $\$k\$$ -th homology embedding and focus on cases reminiscent of spectral clustering. Namely, we analyze the $\{\text{em connected sum}\}$ of manifolds as a perturbation to the direct sum of their homology embeddings. We propose an algorithm to factorize the homology embedding into subspaces corresponding to a manifold's simplest topological components. The proposed framework is applied to the $\{\text{em shortest homologous loop detection}\}$ problem, a problem known to be NP-hard in general. Our spectral loop detection algorithm scales better than existing methods and is effective on diverse data such as point clouds and images.

[Breaking the Moments Condition Barrier: No-Regret Algorithm for Bandits with Super Heavy-Tailed Payoffs](#)

- Han Zhong, Jiayi Huang, Lin Yang, Liwei Wang
- abstract@[open-review](#): Despite a large amount of effort in dealing with heavy-tailed error in machine learning, little is known when moments of the error can become non-existent: the random noise $\$eta\$$ satisfies $\Pr[\eta > ly] \leq 1/y^{\alpha} \text{ for some } \alpha > 0\$$. We make the first attempt to actively handle such super heavy-tailed noise in bandit learning problems: We propose a novel robust statistical estimator, mean of medians, which estimates a random variable by computing the empirical mean of a sequence of empirical medians. We then present a generic reductionist algorithmic framework for solving bandit learning problems (including multi-armed and linear bandit problem): the mean of medians estimator can be applied to nearly any bandit learning algorithm as a black-box filtering for its reward signals and obtain similar regret bound as if the reward is sub-Gaussian. We show that the regret bound is near-optimal even with very heavy-tailed noise. We also empirically demonstrate the effectiveness of the proposed algorithm, which further corroborates our theoretical results.

[A nonparametric method for gradual change problems with statistical guarantees](#)

- Lizhen Nie, Dan Nicolae
- abstract@[open-review](#): We consider the detection and localization of gradual changes in the distribution of a sequence of time-ordered observations. Existing literature focuses mostly on the simpler abrupt setting which assumes a discontinuity jump in distribution, and is unrealistic for some applied settings. We propose a general method for detecting and localizing gradual changes that does not require any specific data generating model, any particular data type, or any prior knowledge about which features of the distribution are subject to change. Despite relaxed assumptions, the proposed method possesses proven theoretical guarantees for both detection and localization.

[Nested Graph Neural Networks](#)

- Muhan Zhang, Pan Li
- abstract@[open-review](#): Graph neural network (GNN)'s success in graph classification is closely related to the Weisfeiler-Lehman (1-WL) algorithm. By iteratively aggregating neighboring node features to a center node, both 1-WL and GNN obtain a node representation that encodes a rooted subtree around the center node. These rooted subtree representations are then pooled into a single representation to represent the whole graph. However, rooted subtrees are of limited expressiveness to represent a non-tree graph. To address it, we propose Nested Graph Neural Networks (NGNNs). NGNN represents a graph with rooted subgraphs instead of rooted subtrees, so that two graphs sharing many identical subgraphs (rather than subtrees) tend to have similar representations. The key is to make each node representation encode a subgraph around it more than a subtree. To achieve this, NGNN extracts a local subgraph around each node and applies a base GNN to each subgraph to learn a subgraph representation. The whole-graph representation is then obtained by pooling these subgraph representations. We provide a rigorous theoretical analysis showing that NGNN is strictly more powerful than 1-WL. In particular, we proved that NGNN can discriminate almost all r-regular graphs, where 1-WL always fails. Moreover, unlike other more powerful GNNs, NGNN only introduces a constant-factor higher time complexity than standard GNNs. NGNN is a plug-and-play framework that can be combined with various base GNNs. We test NGNN with different base GNNs on several benchmark datasets. NGNN uniformly improves their performance and shows highly competitive performance on all datasets.

[Multimodal and Multilingual Embeddings for Large-Scale Speech Mining](#)

- Paul-Ambroise Duquenne, Hongyu Gong, Holger Schwenk
- abstract@[open-review](#): We present an approach to encode a speech signal into a fixed-size representation which minimizes the cosine loss with the existing massively multilingual LASER text embedding space. Sentences are close in this embedding space, independently of their language and modality, either text or audio. Using a similarity metric in that multimodal embedding space, we perform mining of audio in German, French, Spanish and English from LibriVox against billions of sentences from Common Crawl. This yielded more than twenty thousand hours of aligned speech translations. To evaluate the automatically mined speech/text corpora, we train neural speech translation systems for several languages pairs. Adding the mined data, achieves significant improvements in the BLEU score on the CoVoST2 and the MUST-C test sets with respect to a very competitive baseline. Our approach can also be used to directly perform speech-to-speech mining, without the need to first transcribe or translate the data. We obtain more than one thousand three hundred hours of aligned speech in French, German, Spanish and English. This speech corpus has the potential to boost research in speech-to-speech translation which suffers from scarcity of natural end-to-end training data. All the mined multimodal corpora will be made freely available.

[Necessary and sufficient graphical conditions for optimal adjustment sets in causal graphical models with hidden variables](#)

- Jakob Runge
- abstract@[open-review](#): The problem of selecting optimal backdoor adjustment sets to estimate causal effects in graphical models with hidden and conditioned variables is addressed. Previous work has defined optimality as achieving the smallest asymptotic estimation variance and derived an optimal set for the case without hidden variables. For the case with hidden variables there can be settings where no optimal set exists and currently only a sufficient graphical optimality criterion of limited applicability has been derived. In the present work optimality is characterized as maximizing a certain adjustment information which allows to derive a necessary and sufficient graphical criterion for the existence of an optimal adjustment set and a definition and algorithm to construct it. Further, the optimal set is valid if and only if a valid adjustment set exists and has higher (or equal) adjustment information than the Adjust-set proposed in Perković et al. [Journal of Machine Learning Research, 18: 1–62, 2018] for any graph. The results translate to minimal asymptotic estimation variance for a class of estimators whose asymptotic variance follows a certain information-theoretic relation. Numerical experiments indicate that the asymptotic results also hold for relatively small sample sizes and that the optimal adjustment set or minimized variants thereof often yield better variance also beyond that estimator class. Surprisingly, among the randomly created setups more than 90% fulfill the optimality conditions indicating that also in many real-world scenarios graphical optimality may hold.

[On Blame Attribution for Accountable Multi-Agent Sequential Decision Making](#)

- Stelios Triantafyllou, Adish Singla, Goran Radanovic
- abstract@[open-review](#): Blame attribution is one of the key aspects of accountable decision making, as it provides means to quantify the responsibility of an agent for a decision making outcome. In this paper, we study blame attribution in the context of cooperative multi-agent sequential decision making. As a particular setting of interest, we focus on cooperative decision making formalized by Multi-Agent Markov Decision Processes (MMDPs), and we analyze different blame attribution methods derived from or inspired by existing concepts in cooperative game theory. We formalize desirable properties of blame attribution in the setting of interest, and we analyze the relationship between these properties and the studied blame attribution methods. Interestingly, we show that some of the well known blame attribution methods, such as Shapley value, are not performance-incentivizing, while others, such as Banzhaf index, may over-blame agents. To mitigate these value misalignment and fairness issues, we introduce a novel blame attribution method, unique in the set of properties it satisfies, which trade-offs explanatory power (by under-blaming agents) for the aforementioned properties. We further show how to account for uncertainty about agents' decision making policies, and we experimentally: a) validate the qualitative properties of the studied blame attribution methods, and b) analyze their robustness to uncertainty.

[FLEX: Unifying Evaluation for Few-Shot NLP](#)

- Jonathan Bragg, Arman Cohan, Kyle Lo, Iz Beltagy
- abstract@[open-review](#): Few-shot NLP research is highly active, yet conducted in disjoint research threads with evaluation suites that lack challenging-yet-realistic testing setups and fail to employ careful experimental design. Consequently, the community does not know which techniques perform best or even if they outperform simple baselines. In response, we formulate the FLEX Principles, a set of requirements and best practices for unified, rigorous, valid, and cost-sensitive few-shot NLP evaluation. These principles include Sample Size Design, a novel approach to benchmark design that optimizes statistical accuracy and precision while keeping evaluation costs manageable. Following the principles, we release the FLEX benchmark, which includes four few-shot transfer settings, zero-shot evaluation, and a public leaderboard that covers diverse NLP tasks. In addition, we present UniFew, a prompt-based model for few-shot learning that unifies pretraining and finetuning prompt formats, eschewing complex machinery of recent prompt-based approaches in adapting downstream task formats to language model pretraining objectives. We demonstrate that despite simplicity, UniFew achieves results competitive with both popular meta-learning and prompt-based approaches.

[A flow-based latent state generative model of neural population responses to natural images](#)

- Mohammad Bashiri, Edgar Walker, Konstantin-Klemens Lurz, Akshay Jagadish, Taliah Muhammad, Zhiwei Ding, Zhuokun Ding, Andreas Tolias, Fabian Sinz
- abstract@[open-review](#): We present a joint deep neural system identification model for two major sources of neural variability: stimulus-driven and stimulus-conditioned fluctuations. To this end, we combine (1) state-of-the-art deep networks for stimulus-driven activity and (2) a flexible, normalizing flow-based generative model to capture the stimulus-conditioned variability including noise correlations. This allows us to train the model end-to-end without the need for sophisticated probabilistic approximations associated with many latent state models for stimulus-conditioned fluctuations. We train the model on the responses of thousands of neurons from multiple areas of the mouse visual cortex to natural images. We show that our model outperforms previous state-of-the-art models in predicting the distribution of neural population responses to novel stimuli, including shared stimulus-conditioned variability. Furthermore, it successfully learns known latent factors of the population responses that are related to behavioral variables such as pupil dilation, and other factors that vary systematically with brain area or retinotopic location. Overall, our model accurately accounts for two critical sources of neural variability while avoiding several complexities associated with many existing latent state models. It thus provides a useful tool for uncovering the interplay between different factors that contribute to variability in neural activity.

[Learnable Fourier Features for Multi-dimensional Spatial Positional Encoding](#)

- Yang Li, Si Si, Gang Li, Cho-Jui Hsieh, Samy Bengio
- abstract@[open-review](#): Attentional mechanisms are order-invariant. Positional encoding is a crucial component to allow attention-based deep model architectures such as Transformer to address sequences or images where the position of information matters. In this paper, we propose a novel positional encoding method based on learnable Fourier features. Instead of hard-coding each position as a token or a vector, we represent each position, which can be multi-dimensional, as a trainable encoding based on learnable Fourier feature mapping, modulated with a multi-layer perceptron. The representation is particularly advantageous for a spatial multi-dimensional position, e.g., pixel positions on an image, where $\$L_2\$$ distances or more complex positional relationships need to be captured. Our experiments based on several public benchmark tasks show that our learnable Fourier feature representation for multi-dimensional positional encoding outperforms existing methods by both improving the accuracy and allowing faster convergence.

[Doubly Robust Thompson Sampling with Linear Payoffs](#)

- Wonyoung Kim, Gi-Soo Kim, Myunghee Cho Paik
- abstract@[open-review](#): A challenging aspect of the bandit problem is that a stochastic reward is observed only for the chosen arm and the rewards of other arms remain missing. The dependence of the arm choice on the past context and reward pairs compounds the complexity of regret analysis. We propose a novel multi-armed contextual bandit algorithm called Doubly Robust Thompson Sampling (DRTS) employing the doubly-robust estimator used in missing data literature to Thompson Sampling with contexts (\texttt{LinTS}). Different from previous works relying on missing data techniques (Dimakopoulou et al. [2019], Kim and Paik [2019]), the proposed algorithm is designed to allow a novel additive regret decomposition leading to an improved regret bound with the order of $\$O(\phi^{-2}\sqrt{T})$, where ϕ is the minimum eigenvalue of the covariance matrix of contexts. This is the first regret bound of \texttt{LinTS} using ϕ without d , where d is the dimension of the context. Applying the relationship between ϕ and d , the regret bound of the proposed algorithm is $\tilde{O}(d\sqrt{T})$ in many practical scenarios, improving the bound of \texttt{LinTS} by a factor of \sqrt{d} . A benefit of the proposed method is that it uses all the context data, chosen or not chosen, thus allowing to circumvent the technical definition of unsaturated arms used in theoretical analysis of \texttt{LinTS}. Empirical studies show the advantage of the proposed algorithm over \texttt{LinTS}.

[A Computationally Efficient Method for Learning Exponential Family Distributions](#)

- Abhin Shah, Devavrat Shah, Gregory Wornell
- abstract@[open-review](#): We consider the question of learning the natural parameters of a k parameter \textit{minimal} exponential family from i.i.d. samples in a computationally and statistically efficient manner. We focus on the setting where the support as well as the natural parameters are appropriately bounded. While the traditional maximum likelihood estimator for this class of exponential family is consistent, asymptotically normal, and asymptotically efficient, evaluating it is computationally hard. In this work, we propose a computationally efficient estimator that is consistent as well as asymptotically normal under mild conditions. We provide finite sample guarantees to achieve an (λ) error of α in the parameter estimation with sample complexity $O(\mathrm{poly}(k/\alpha))$ and computational complexity $O(\mathrm{poly}(k/\alpha))$. To establish these results, we show that, at the population level, our method can be viewed as the maximum likelihood estimation of a re-parameterized distribution belonging to the same class of exponential family.

[Rethinking Neural Operations for Diverse Tasks](#)

- Nicholas Roberts, Mikhail Khodak, Tri Dao, Liam Li, Christopher Ré, Ameet Talwalkar
- abstract@[open-review](#): An important goal of AutoML is to automate-away the design of neural networks on new tasks in under-explored domains. Motivated by this goal, we study the problem of enabling users to discover the right neural operations given data from their specific domain. We introduce a search space of operations called XD-Operations that mimic the inductive bias of standard multi-channel convolutions while being much more expressive: we prove that it includes many named operations across multiple application areas. Starting with any standard backbone such as ResNet, we show how to transform it into a search space over XD-operations and how to traverse the space using a simple weight sharing scheme. On a diverse set of tasks—solving PDEs, distance prediction for protein folding, and music modeling—our approach consistently yields models with lower error than baseline networks and often even lower error than expert-designed domain-specific approaches.

[Motif-based Graph Self-Supervised Learning for Molecular Property Prediction](#)

- ZAXI ZHANG, Qi Liu, Hao Wang, Chengqiang Lu, Chee-Kong Lee
- abstract@[open-review](#): Predicting molecular properties with data-driven methods has drawn much attention in recent years. Particularly, Graph Neural Networks (GNNs) have demonstrated remarkable success in various molecular generation and prediction tasks. In cases where labeled data is scarce, GNNs can be pre-trained on unlabeled molecular data to first learn the general semantic and structural information before being finetuned for specific tasks. However, most existing self-supervised pretraining frameworks for GNNs only focus on node-level or graph-level tasks. These approaches cannot capture the rich information in subgraphs or graph motifs. For example, functional groups (frequently-occurred subgraphs in molecular graphs) often carry indicative information about the molecular properties. To bridge this gap, we propose Motif-based Graph Self-supervised Learning (MGSSL) by introducing a novel self-supervised motif generation framework for GNNs. First, for motif extraction from molecular graphs, we design a molecule fragmentation method that leverages a retrosynthesis-based algorithm BRICS and additional rules for controlling the size of motif vocabulary. Second, we design a general motif-based generative pretraining framework in which GNNs are asked to make topological and label predictions. This generative framework can be implemented in two different ways, i.e., breadth-first or depth-first. Finally, to take the multi-scale information in molecular graphs into consideration, we introduce a multi-level self-supervised pre-training. Extensive experiments on various downstream benchmark tasks show that our methods outperform all state-of-the-art baselines.

[On Inductive Biases for Heterogeneous Treatment Effect Estimation](#)

- Alicia Curth, Mihaela van der Schaar

- abstract@[open-review](#): We investigate how to exploit structural similarities of an individual's potential outcomes (POs) under different treatments to obtain better estimates of conditional average treatment effects in finite samples. Especially when it is unknown whether a treatment has an effect at all, it is natural to hypothesize that the POs are similar -- yet, some existing strategies for treatment effect estimation employ regularization schemes that implicitly encourage heterogeneity even when it does not exist and fail to fully make use of shared structure. In this paper, we investigate and compare three end-to-end learning strategies to overcome this problem -- based on regularization, reparametrization and a flexible multi-task architecture -- each encoding inductive bias favoring shared behavior across POs. To build understanding of their relative strengths, we implement all strategies using neural networks and conduct a wide range of semi-synthetic experiments. We observe that all three approaches can lead to substantial improvements upon numerous baselines and gain insight into performance differences across various experimental settings.

[DP-SSL: Towards Robust Semi-supervised Learning with A Few Labeled Samples](#)

- Yi Xu, Jiandong Ding, Lu Zhang, Shuigeng Zhou
- abstract@[open-review](#): The scarcity of labeled data is a critical obstacle to deep learning. Semi-supervised learning (SSL) provides a promising way to leverage unlabeled data by pseudo labels. However, when the size of labeled data is very small (say a few labeled samples per class), SSL performs poorly and unstably, possibly due to the low quality of learned pseudo labels. In this paper, we propose a new SSL method called DP-SSL that adopts an innovative data programming (DP) scheme to generate probabilistic labels for unlabeled data. Different from existing DP methods that rely on human experts to provide initial labeling functions (LFs), we develop a multiple-choice learning~(MCL) based approach to automatically generate LFs from scratch in SSL style. With the noisy labels produced by the LFs, we design a label model to resolve the conflict and overlap among the noisy labels, and finally infer probabilistic labels for unlabeled samples. Extensive experiments on four standard SSL benchmarks show that DP-SSL can provide reliable labels for unlabeled data and achieve better classification performance on test sets than existing SSL methods, especially when only a small number of labeled samples are available. Concretely, for CIFAR-10 with only 40 labeled samples, DP-SSL achieves 93.82% annotation accuracy on unlabeled data and 93.46% classification accuracy on test data, which are higher than the SOTA results.

[Transformer in Transformer](#)

- Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing XU, Yunhe Wang
- abstract@[open-review](#): Transformer is a new kind of neural architecture which encodes the input data as powerful features via the attention mechanism. Basically, the visual transformers first divide the input images into several local patches and then calculate both representations and their relationship. Since natural images are of high complexity with abundant detail and color information, the granularity of the patch dividing is not fine enough for excavating features of objects in different scales and locations. In this paper, we point out that the attention inside these local patches are also essential for building visual transformers with high performance and we explore a new architecture, namely, Transformer iN Transformer (TNT). Specifically, we regard the local patches (\eg, 16\$\times\$16) as "visual sentences" and present to further divide them into smaller patches (\eg, 4\$\times\$4) as "visual words". The attention of each word will be calculated with other words in the given visual sentence with negligible computational costs. Features of both words and sentences will be aggregated to enhance the representation ability. Experiments on several benchmarks demonstrate the effectiveness of the proposed TNT architecture, \eg, we achieve an 81.5\% top-1 accuracy on the ImageNet, which is about 1.7\% higher than that of the state-of-the-art visual transformer with similar computational cost. The PyTorch code is available at \url{https://github.com/huawei-noah/CV-Backbones}, and the MindSpore code is available at \url{https://gitee.com/mindspore/models/tree/master/research/cv/TNT}.

[Adversarial Graph Augmentation to Improve Graph Contrastive Learning](#)

- Susheel Suresh, Pan Li, Cong Hao, Jennifer Neville
- abstract@[open-review](#): Self-supervised learning of graph neural networks (GNN) is in great need because of the widespread label scarcity issue in real-world graph/network data. Graph contrastive learning (GCL), by training GNNs to maximize the correspondence between the representations of the same graph in its different augmented forms, may yield robust and transferable GNNs even without using labels. However, GNNs trained by traditional GCL often risk capturing redundant graph features and thus may be brittle and provide sub-par performance in downstream tasks. Here, we propose a novel principle, termed adversarial-GCL (\textit{AD-GCL}), which enables GNNs to avoid capturing redundant information during the training by optimizing adversarial graph augmentation strategies used in GCL. We pair AD-GCL with theoretical explanations and design a practical instantiation based on trainable edge-dropping graph augmentation. We experimentally validate AD-GCL by comparing with the state-of-the-art GCL methods and achieve performance gains of up-to~14\% in unsupervised, ~6\% in transfer and~3\% in semi-supervised learning settings overall with 18 different benchmark datasets for the tasks of molecule property regression and classification, and social network classification.

[Online Control of Unknown Time-Varying Dynamical Systems](#)

- Edgar Minasyan, Paula Gradu, Max Simchowitz, Elad Hazan
- abstract@[open-review](#): We study online control of time-varying linear systems with unknown dynamics in the nonstochastic control model. At a high level, we demonstrate that this setting is \emph{qualitatively harder} than that of either unknown time-invariant or known time-varying dynamics, and complement our negative results with algorithmic upper bounds in regimes where sublinear regret is possible. More specifically, we study regret bounds with respect to common classes of policies: Disturbance Action (SLS), Disturbance Response (Youla), and linear feedback policies. While these three classes are essentially equivalent for LTI systems, we demonstrate that these equivalences break down for time-varying systems. We prove a lower bound that no algorithm can obtain sublinear regret with respect to the first two classes unless a certain measure of system variability also scales sublinearly in the horizon. Furthermore, we show that offline planning over the state linear feedback policies is NP-hard, suggesting hardness of the online learning problem. On the positive side, we give an efficient algorithm that attains a sublinear regret bound against the class of Disturbance Response policies up to the aforementioned system variability term. In fact, our algorithm enjoys sublinear \emph{adaptive} regret bounds, which is a strictly stronger metric than standard regret and is more appropriate for time-varying systems. We sketch extensions to Disturbance Action policies and partial observation, and propose an inefficient algorithm for regret against linear state feedback policies.

[Contrastive Reinforcement Learning of Symbolic Reasoning Domains](#)

- Gabriel Poesia, WenXin Dong, Noah Goodman
- abstract@[open-review](#): Abstract symbolic reasoning, as required in domains such as mathematics and logic, is a key component of human intelligence. Solvers for these domains have important applications, especially to computer-assisted education. But learning to solve symbolic problems is challenging for machine learning algorithms. Existing models either learn from human solutions or use hand-engineered features, making them expensive to apply in new domains. In this paper, we instead consider symbolic domains as simple environments where states and actions are given as unstructured text, and binary rewards indicate whether a problem is solved. This flexible setup makes it easy to specify new domains, but search and planning become challenging. We introduce five environments inspired by the Mathematics Common Core Curriculum, and observe that existing Reinforcement Learning baselines perform poorly. We then present a novel learning algorithm, Contrastive Policy Learning (ConPoLe) that explicitly optimizes the InfoNCE loss, which lower bounds the mutual information between the current state and next states that continue on a path to the solution. ConPoLe successfully solves all four domains. Moreover, problem representations learned by ConPoLe enable accurate prediction of the categories of problems in a real mathematics curriculum. Our results suggest new directions for reinforcement learning in symbolic domains, as well as applications to mathematics education.

[Spatial Ensemble: a Novel Model Smoothing Mechanism for Student-Teacher Framework](#)

- Tengteng Huang, Yifan Sun, Xun Wang, Haotian Yao, Chi Zhang
- abstract@[open-review](#): Model smoothing is of central importance for obtaining a reliable teacher model in the student-teacher framework, where the teacher generates surrogate supervision signals to train the student. A popular model smoothing method is the Temporal Moving Average (TMA), which continuously averages the teacher parameters with the up-to-date student parameters. In this paper, we propose "Spatial Ensemble", a novel model smoothing mechanism in parallel with TMA. Spatial Ensemble randomly picks up a small fragment of the student model to directly replace the corresponding fragment of the teacher model. Consequentially, it stitches different fragments of historical student models into a unity, yielding the "Spatial Ensemble" effect. Spatial Ensemble obtains comparable student-teacher learning performance by itself and demonstrates valuable complementarity with temporal moving average. Their integration, named Spatial-Temporal Smoothing, brings general (sometimes significant) improvement to the student-teacher learning framework on a variety of state-of-the-art methods. For example, based on the self-supervised method BYOL, it yields +0.9% top-1 accuracy improvement on ImageNet, while based on the semi-supervised approach FixMatch, it increases the top-1 accuracy by around +6% on CIFAR-10 when only few training labels are available. Codes and models are available at: https://github.com/tengteng95/Spatial_Ensemble.

[Probabilistic Tensor Decomposition of Neural Population Spiking Activity](#)

- Hugo Soulat, Sepiedeh Keshavarzi, Troy Margrie, Maneesh Sahani
- abstract@[open-review](#): The firing of neural populations is coordinated across cells, in time, and across experimental conditions or repeated experimental trials; and so a full understanding of the computational significance of neural responses must be based on a separation of these different contributions to structured activity. Tensor decomposition is an approach to untangling the influence of multiple factors in data that is common in many fields. However, despite some recent interest in neuroscience, wider applicability of the approach is hampered by the lack of a full probabilistic treatment allowing principled inference of a decomposition from non-Gaussian spike-count data. Here, we extend the Pólya-Gamma (PG) augmentation, previously used in sampling-based Bayesian inference, to implement scalable variational inference in non-conjugate spike-count models. Using this new approach, we develop techniques related to automatic relevance determination to infer the most appropriate tensor rank, as well as to incorporate priors based on known brain anatomy such as the segregation of cell response properties by brain area. We apply the model to neural recordings taken under conditions of visual-vestibular sensory integration, revealing how the encoding of self- and visual-motion signals is modulated by the sensory information available to the animal.

[Recurrent Bayesian Classifier Chains for Exact Multi-Label Classification](#)

- Walter Gerych, Tom Hartvigsen, Luke Buquicchio, Emmanuel Agu, Elke A. Rundensteiner
- abstract@[open-review](#): Exact multi-label classification is the task of assigning each datapoint a set of class labels such that the assigned set exactly matches the ground truth. Optimizing for exact multi-label classification is important in domains where missing a single label can be especially costly, such as in object detection for autonomous vehicles or symptom classification for disease diagnosis. Recurrent Classifier Chains (RCCs), a recurrent neural network extension of ensemble-based classifier chains, are the state-of-the-art exact multi-label classification method for maximizing subset accuracy. However, RCCs iteratively predict classes with an unprincipled ordering, and therefore indiscriminately condition class probabilities. These disadvantages make RCCs prone to predicting inaccurate label sets. In this work we propose Recurrent Bayesian Classifier Chains (RBCCs), which learn a Bayesian network of class dependencies and leverage this network in order to condition the prediction of child nodes only on their parents. By conditioning predictions in this way, we perform principled and non-noisy class prediction. We demonstrate the effectiveness of our RBCC method on a variety of real-world multi-label datasets, where we routinely outperform the state of the art methods for exact multi-label classification.

[Wasserstein Flow Meets Replicator Dynamics: A Mean-Field Analysis of Representation Learning in Actor-Critic](#)

- Yufeng Zhang, Siyu Chen, Zhuoran Yang, Michael Jordan, Zhaoran Wang
- abstract@[open-review](#): Actor-critic (AC) algorithms, empowered by neural networks, have had significant empirical success in recent years. However, most of the existing theoretical support for AC algorithms focuses on the case of linear function approximations, or linearized neural networks, where the feature representation is fixed throughout training. Such a limitation fails to capture the key aspect of representation learning in neural AC, which is pivotal in practical problems. In this work, we take a mean-field perspective on the evolution and convergence of feature-based neural AC. Specifically, we consider a version of AC where the actor and critic are represented by overparameterized two-layer neural networks and are updated with two-timescale learning rates. The critic is updated by temporal-difference (TD) learning with a larger stepsize while the actor is updated via proximal policy optimization (PPO) with a smaller stepsize. In the continuous-time and infinite-width limiting regime, when the timescales are properly separated, we prove that neural AC finds the globally optimal policy at a sublinear rate. Additionally, we prove that the feature representation induced by the critic network is allowed to evolve within a neighborhood of the initial one.

[Assessing Fairness in the Presence of Missing Data](#)

- Yiliang Zhang, Qi Long
- abstract@[open-review](#): Missing data are prevalent and present daunting challenges in real data analysis. While there is a growing body of literature on fairness in analysis of fully observed data, there has been little theoretical work on investigating fairness in analysis of incomplete data. In practice, a popular analytical approach for dealing with missing data is to use only the set of complete cases, i.e., observations with all features fully observed to train a prediction algorithm. However, depending on the missing data mechanism, the distribution of complete cases and the distribution of the complete data may be substantially different. When the goal is to develop a fair algorithm in the complete data domain where there are no missing values, an algorithm that is fair in the complete case domain may show disproportionate bias towards some marginalized groups in the complete data domain. To fill this significant gap, we study the problem of estimating fairness in the complete data domain for an arbitrary model evaluated merely using complete cases. We provide upper and lower bounds on the fairness estimation error and conduct numerical experiments to assess our theoretical results. Our work provides the first known theoretical results on fairness guarantee in analysis of incomplete data.

[Adversarial Attack Generation Empowered by Min-Max Optimization](#)

- Jingkang Wang, Tianyun Zhang, Sijia Liu, Pin-Yu Chen, Jiacen Xu, Makan Fardad, Bo Li
- abstract@[open-review](#): The worst-case training principle that minimizes the maximal adversarial loss, also known as adversarial training (AT), has shown to be a state-of-the-art approach for enhancing adversarial robustness. Nevertheless, min-max optimization beyond the purpose of AT has not been rigorously explored in the adversarial context. In this paper, we show how a general notion of min-max optimization over multiple domains can be leveraged to the design of different types of adversarial attacks. In particular, given a set of risk sources, minimizing the worst-case attack loss can be reformulated as a min-max problem by introducing domain weights that are maximized over the probability simplex of the domain set. We showcase this unified framework in three attack generation problems -- attacking model ensembles, devising universal perturbation under multiple inputs, and crafting attacks resilient to data transformations. Extensive experiments demonstrate that our approach leads to substantial attack improvement over the existing heuristic strategies as well as robustness improvement over state-of-the-art defense methods against multiple perturbation types. Furthermore, we find that the self-adjusted domain weights learned from min-max optimization can provide a holistic tool to explain the difficulty level of attack across domains.

[Safe Pontryagin Differentiable Programming](#)

- Wanxin Jin, Shaoshuai Mou, George J. Pappas

- abstract@[open-review](#): We propose a Safe Pontryagin Differentiable Programming (Safe PDP) methodology, which establishes a theoretical and algorithmic framework to solve a broad class of safety-critical learning and control tasks---problems that require the guarantee of safety constraint satisfaction at any stage of the learning and control progress. In the spirit of interior-point methods, Safe PDP handles different types of system constraints on states and inputs by incorporating them into the cost or loss through barrier functions. We prove three fundamentals of the proposed Safe PDP: first, both the solution and its gradient in the backward pass can be approximated by solving their more efficient unconstrained counterparts; second, the approximation for both the solution and its gradient can be controlled for arbitrary accuracy by a barrier parameter; and third, importantly, all intermediate results throughout the approximation and optimization strictly respect the constraints, thus guaranteeing safety throughout the entire learning and control process. We demonstrate the capabilities of Safe PDP in solving various safety-critical tasks, including safe policy optimization, safe motion planning, and learning MPCs from demonstrations, on different challenging systems such as 6-DoF maneuvering quadrotor and 6-DoF rocket powered landing.

[Class-Disentanglement and Applications in Adversarial Detection and Defense](#)

- Kaiwen Yang, Tianyi Zhou, Yonggang Zhang, Xinmei Tian, Dacheng Tao
- abstract@[open-review](#): What is the minimum necessary information required by a neural net $D(\cdot)$ from an image x to accurately predict its class? Extracting such information in the input space from x can allocate the areas $D(\cdot)$ mainly attending to and shed novel insights to the detection and defense of adversarial attacks. In this paper, we propose "class-disentanglement" that trains a variational autoencoder $G(\cdot)$ to extract this class-dependent information as $x - G(x)$ via a trade-off between reconstructing x by $G(x)$ and classifying x by $D(x - G(x))$, where the former competes with the latter in decomposing x so the latter retains only necessary information for classification in $x - G(x)$. We apply it to both clean images and their adversarial images and discover that the perturbations generated by adversarial attacks mainly lie in the class-dependent part $x - G(x)$. The decomposition results also provide novel interpretations to classification and attack models. Inspired by these observations, we propose to conduct adversarial detection and adversarial defense respectively on $x - G(x)$ and $G(x)$, which consistently outperform the results on the original x . In experiments, this simple approach substantially improves the detection and defense against different types of adversarial attacks.

[Active 3D Shape Reconstruction from Vision and Touch](#)

- Edward Smith, David Meger, Luis Pineda, Roberto Calandra, Jitendra Malik, Adriana Romero Soriano, Michal Drozdzał
- abstract@[open-review](#): Humans build 3D understandings of the world through active object exploration, using jointly their senses of vision and touch. However, in 3D shape reconstruction, most recent progress has relied on static datasets of limited sensory data such as RGB images, depth maps or haptic readings, leaving the active exploration of the shape largely unexplored. In active touch sensing for 3D reconstruction, the goal is to actively select the tactile readings that maximize the improvement in shape reconstruction accuracy. However, the development of deep learning-based active touch models is largely limited by the lack of frameworks for shape exploration. In this paper, we focus on this problem and introduce a system composed of: 1) a haptic simulator leveraging high spatial resolution vision-based tactile sensors for active touching of 3D objects; 2) a mesh-based 3D shape reconstruction model that relies on tactile or visuotactile signals; and 3) a set of data-driven solutions with either tactile or visuotactile priors to guide the shape exploration. Our framework enables the development of the first fully data-driven solutions to active touch on top of learned models for object understanding. Our experiments show the benefits of such solutions in the task of 3D shape understanding where our models consistently outperform natural baselines. We provide our framework as a tool to foster future research in this direction.

[CAPE: Encoding Relative Positions with Continuous Augmented Positional Embeddings](#)

- Tatiana Likhomanenko, Qiantong Xu, Gabriel Synnaeve, Ronan Collobert, Alex Rogozhnikov
- abstract@[open-review](#): Without positional information, attention-based Transformer neural networks are permutation-invariant. Absolute or relative positional embeddings are the most popular ways to feed Transformer models with positional information. Absolute positional embeddings are simple to implement, but suffer from generalization issues when evaluating on sequences longer than seen at training time. Relative positions are more robust to input length change, but are more complex to implement and yield inferior model throughput due to extra computational and memory costs. In this paper, we propose an augmentation-based approach (CAPE) for absolute positional embeddings, which keeps the advantages of both absolute (simplicity and speed) and relative positional embeddings (better generalization). In addition, our empirical evaluation on state-of-the-art models in machine translation, image and speech recognition demonstrates that CAPE leads to better generalization performance as well as increased stability with respect to training hyper-parameters.

[Multi-armed Bandit Requiring Monotone Arm Sequences](#)

- Ningyan Chen
- abstract@[open-review](#): In many online learning or multi-armed bandit problems, the taken actions or pulled arms are ordinal and required to be monotone over time. Examples include dynamic pricing, in which the firms use markup pricing policies to please early adopters and deter strategic waiting, and clinical trials, in which the dose allocation usually follows the dose escalation principle to prevent dose limiting toxicities. We consider the continuum-armed bandit problem when the arm sequence is required to be monotone. We show that when the unknown objective function is Lipschitz continuous, the regret is $O(T)$. When in addition the objective function is unimodal or quasiconcave, the regret is $\tilde{O}(T^{3/4})$ under the proposed algorithm, which is also shown to be the optimal rate. This deviates from the optimal rate $\tilde{O}(T^{2/3})$ in the continuous-armed bandit literature and demonstrates the cost to the learning efficiency brought by the monotonicity requirement.

[Gradient Driven Rewards to Guarantee Fairness in Collaborative Machine Learning](#)

- Xinyi Xu, Lingjuan Lyu, Xingjun Ma, Chenglin Miao, Chuan Sheng Foo, Bryan Kian Hsiang Low
- abstract@[open-review](#): In collaborative machine learning(CML), multiple agents pool their resources(e.g., data) together for a common learning task. In realistic CML settings where the agents are self-interested and not altruistic, they may be unwilling to share data or model information without adequate rewards. Furthermore, as the data/model information shared by the agents may differ in quality, designing rewards which are fair to them is important so that they would not feel exploited nor discouraged from sharing. In this paper, we adopt federated learning as the CML paradigm, propose a novel cosine gradient Shapley value(CGSV) to fairly evaluate the expected marginal contribution of each agent's uploaded model parameter update/gradient without needing an auxiliary validation dataset, and based on the CGSV, design a novel training-time gradient reward mechanism with a fairness guarantee by sparsifying the aggregated parameter update/gradient downloaded from the server as reward to each agent such that its resulting quality is commensurate to that of the agent's uploaded parameter update/gradient. We empirically demonstrate the effectiveness of our fair gradient reward mechanism on multiple benchmark datasets in terms of fairness, predictive performance, and time overhead.

[Generalizable Imitation Learning from Observation via Inferring Goal Proximity](#)

- Youngwoon Lee, Andrew Szot, Shao-Hua Sun, Joseph J. Lim
- abstract@[open-review](#): Task progress is intuitive and readily available task information that can guide an agent closer to the desired goal. Furthermore, a task progress estimator can generalize to new situations. From this intuition, we propose a simple yet effective imitation learning from observation method for a goal-directed task using a learned goal proximity function as a task progress estimator for better generalization to unseen states and goals. We obtain this goal proximity function from expert demonstrations and online agent experience, and then use the learned goal proximity as a dense reward for policy training. We demonstrate that our proposed method can robustly generalize compared to prior imitation learning methods on a set of goal-directed tasks in navigation, locomotion, and robotic manipulation, even with demonstrations that cover only a part of the states.

[DualNet: Continual Learning, Fast and Slow](#)

- Quang Pham, Chenghao Liu, Steven Hoi
- abstract@[open-review](#): According to Complementary Learning Systems (CLS) theory~\cite{mcclelland1995there} in neuroscience, humans do effective continual learning through two complementary systems: a fast learning system centered on the hippocampus for rapid learning of the specifics and individual experiences, and a slow learning system located in the neocortex for the gradual acquisition of structured knowledge about the environment. Motivated by this theory, we propose a novel continual learning framework named ``DualNet'', which comprises a fast learning system for supervised learning of pattern-separated representation from specific tasks and a slow learning system for unsupervised representation learning of task-agnostic general representation via a Self-Supervised Learning (SSL) technique. The two fast and slow learning systems are complementary and work seamlessly in a holistic continual learning framework. Our extensive experiments on two challenging continual learning benchmarks of CORE50 and miniImageNet show that DualNet outperforms state-of-the-art continual learning methods by a large margin. We further conduct ablation studies of different SSL objectives to validate DualNet's efficacy, robustness, and scalability. Code is publicly available at \url{https://github.com/phquang/DualNet}.

[Deformable Butterfly: A Highly Structured and Sparse Linear Transform](#)

- Rui Lin, Jie Ran, King Hung Chiu, Graziano Chesi, Ngai Wong
- abstract@[open-review](#): We introduce a new kind of linear transform named Deformable Butterfly (DeBut) that generalizes the conventional butterfly matrices and can be adapted to various input-output dimensions. It inherits the fine-to-coarse-grained learnable hierarchy of traditional butterflies and when deployed to neural networks, the prominent structures and sparsity in a DeBut layer constitutes a new way for network compression. We apply DeBut as a drop-in replacement of standard fully connected and convolutional layers, and demonstrate its superiority in homogenizing a neural network and rendering it favorable properties such as light weight and low inference complexity, without compromising accuracy. The natural complexity-accuracy tradeoff arising from the myriad deformations of a DeBut layer also opens up new rooms for analytical and practical research. The codes and Appendix are publicly available at: <https://github.com/ruilin0212/DeBut>.

[Why Do Pretrained Language Models Help in Downstream Tasks? An Analysis of Head and Prompt Tuning](#)

- Colin Wei, Sang Michael Xie, Tengyu Ma
- abstract@[open-review](#): Pretrained language models have achieved state-of-the-art performance when adapted to a downstream NLP task. However, theoretical analysis of these models is scarce and challenging since the pretraining and downstream tasks can be very different. We propose an analysis framework that links the pretraining and downstream tasks with an underlying latent variable generative model of text -- the downstream classifier must recover a function of the posterior distribution over the latent variables. We analyze head tuning (learning a classifier on top of the frozen pretrained model) and prompt tuning in this setting. The generative model in our analysis is either a Hidden Markov Model (HMM) or an HMM augmented with a latent memory component, motivated by long-term dependencies in natural language. We show that 1) under certain non-degeneracy conditions on the HMM, simple classification heads can solve the downstream task, 2) prompt tuning obtains downstream guarantees with weaker non-degeneracy conditions, and 3) our recovery guarantees for the memory-augmented HMM are stronger than for the vanilla HMM because task-relevant information is easier to recover from the long-term memory. Experiments on synthetically generated data from HMMs back our theoretical findings.

[Learning Diverse Policies in MOBA Games via Macro-Goals](#)

- Yiming Gao, Bei Shi, Xueying Du, Liang Wang, Guangwei Chen, Zhenjie Lian, Fuhao Qiu, GUOAN HAN, Weixuan Wang, Deheng Ye, Qiang Fu, Wei Yang, Lanxiao Huang
- abstract@[open-review](#): Recently, many researchers have made successful progress in building the AI systems for MOBA-game-playing with deep reinforcement learning, such as on Dota 2 and Honor of Kings. Even though these AI systems have achieved or even exceeded human-level performance, they still suffer from the lack of policy diversity. In this paper, we propose a novel Macro-Goals Guided framework, called MGG, to learn diverse policies in MOBA games. MGG abstracts strategies as macro-goals from human demonstrations and trains a Meta-Controller to predict these macro-goals. To enhance policy diversity, MGG samples macro-goals from the Meta-Controller prediction and guides the training process towards these goals. Experimental results on the typical MOBA game Honor of Kings demonstrate that MGG can execute diverse policies in different matches and lineups, and also outperform the state-of-the-art methods over 102 heroes.

[Evaluation of Human-AI Teams for Learned and Rule-Based Agents in Hanabi](#)

- Ho Chit Siu, Jaime Peña, Edenna Chen, Yutai Zhou, Victor Lopez, Kyle Palko, Kimberlee Chang, Ross Allen
- abstract@[open-review](#): Deep reinforcement learning has generated superhuman AI in competitive games such as Go and StarCraft. Can similar learning techniques create a superior AI teammate for human-machine collaborative games? Will humans prefer AI teammates that improve objective team performance or those that improve subjective metrics of trust? In this study, we perform a single-blind evaluation of teams of humans and AI agents in the cooperative card game Hanabi, with both rule-based and learning-based agents. In addition to the game score, used as an objective metric of the human-AI team performance, we also quantify subjective measures of the human's perceived performance, teamwork, interpretability, trust, and overall preference of AI teammate. We find that humans have a clear preference toward a rule-based AI teammate (SmartBot) over a state-of-the-art learning-based AI teammate (Other-Play) across nearly all subjective metrics, and generally view the learning-based agent negatively, despite no statistical difference in the game score. This result has implications for future AI design and reinforcement learning benchmarking, highlighting the need to incorporate subjective metrics of human-AI teaming rather than a singular focus on objective task performance.

[Counterfactual Invariance to Spurious Correlations in Text Classification](#)

- Victor Veitch, Alexander D'Amour, Steve Yadlowsky, Jacob Eisenstein
- abstract@[open-review](#): Informally, a 'spurious correlation' is the dependence of a model on some aspect of the input data that an analyst thinks shouldn't matter. In machine learning, these have a know-it-when-you-see-it character; e.g., changing the gender of a sentence's subject changes a sentiment predictor's output. To check for spurious correlations, we can 'stress test' models by perturbing irrelevant parts of input data and seeing if model predictions change. In this paper, we study stress testing using the tools of causal inference. We introduce counterfactual invariance as a formalization of the requirement that changing irrelevant parts of the input shouldn't change model predictions. We connect counterfactual invariance to out-of-domain model performance, and provide practical schemes for learning (approximately) counterfactual invariant predictors (without access to counterfactual examples). It turns out that both the means and implications of counterfactual invariance depend fundamentally on the true underlying causal structure of the data--in particular, whether the label causes the features or the features cause the label. Distinct causal structures require distinct regularization schemes to induce counterfactual invariance. Similarly, counterfactual invariance implies different domain shift guarantees depending on the underlying causal structure. This theory is supported by empirical results on text classification.

[Better Safe Than Sorry: Preventing Delusive Adversaries with Adversarial Training](#)

- Lue Tao, Lei Feng, Jinfeng Yi, Sheng-Jun Huang, Songcan Chen

- abstract@[open-review](#): Delusive attacks aim to substantially deteriorate the test accuracy of the learning model by slightly perturbing the features of correctly labeled training examples. By formalizing this malicious attack as finding the worst-case training data within a specific ∞ -Wasserstein ball, we show that minimizing adversarial risk on the perturbed data is equivalent to optimizing an upper bound of natural risk on the original data. This implies that adversarial training can serve as a principled defense against delusive attacks. Thus, the test accuracy decreased by delusive attacks can be largely recovered by adversarial training. To further understand the internal mechanism of the defense, we disclose that adversarial training can resist the delusive perturbations by preventing the learner from overly relying on non-robust features in a natural setting. Finally, we complement our theoretical findings with a set of experiments on popular benchmark datasets, which show that the defense withstands six different practical attacks. Both theoretical and empirical results vote for adversarial training when confronted with delusive adversaries.

[Determinantal point processes based on orthogonal polynomials for sampling minibatches in SGD](#)

- Rémi Bardenet, Subhroshekhar Ghosh, Meixia LIN
- abstract@[open-review](#): Stochastic gradient descent (SGD) is a cornerstone of machine learning. When the number N of data items is large, SGD relies on constructing an unbiased estimator of the gradient of the empirical risk using a small subset of the original dataset, called a minibatch. Default minibatch construction involves uniformly sampling a subset of the desired size, but alternatives have been explored for variance reduction. In particular, experimental evidence suggests drawing minibatches from determinantal point processes (DPPs), tractable distributions over minibatches that favour diversity among selected items. However, like in recent work on DPPs for coresets, providing a systematic and principled understanding of how and why DPPs help has been difficult. In this work, we contribute an orthogonal polynomial-based determinantal point process paradigm for performing minibatch sampling in SGD. Our approach leverages the specific data distribution at hand, which endows it with greater sensitivity and power over existing data-agnostic methods. We substantiate our method via a detailed theoretical analysis of its convergence properties, interweaving between the discrete data set and the underlying continuous domain. In particular, we show how specific DPPs and a string of controlled approximations can lead to gradient estimators with a variance that decays faster with the batchsize than under uniform sampling. Coupled with existing finite-time guarantees for SGD on convex objectives, this entails that, for a large enough batchsize and a fixed budget of item-level gradients to evaluate, DPP minibatches lead to a smaller bound on the mean square approximation error than uniform minibatches. Moreover, our estimators are amenable to a recent algorithm that directly samples linear statistics of DPPs (i.e., the gradient estimator) without sampling the underlying DPP (i.e., the minibatch), thereby reducing computational overhead. We provide detailed synthetic as well as real data experiments to substantiate our theoretical claims.

[Revisiting Contrastive Methods for Unsupervised Learning of Visual Representations](#)

- Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Luc V Gool
- abstract@[open-review](#): Contrastive self-supervised learning has outperformed supervised pretraining on many downstream tasks like segmentation and object detection. However, current methods are still primarily applied to curated datasets like ImageNet. In this paper, we first study how biases in the dataset affect existing methods. Our results show that an approach like MoCo works surprisingly well across: (i) object- versus scene-centric, (ii) uniform versus long-tailed and (iii) general versus domain-specific datasets. Second, given the generality of the approach, we try to realize further gains with minor modifications. We show that learning additional invariances - through the use of multi-scale cropping, stronger augmentations and nearest neighbors - improves the representations. Finally, we observe that MoCo learns spatially structured representations when trained with a multi-crop strategy. The representations can be used for semantic segment retrieval and video instance segmentation without finetuning. Moreover, the results are on par with specialized models. We hope this work will serve as a useful study for other researchers.

[Neural Analysis and Synthesis: Reconstructing Speech from Self-Supervised Representations](#)

- Hyeong-Seok Choi, Juheon Lee, Wansoo Kim, Jie Lee, Hoon Heo, Kyogu Lee
- abstract@[open-review](#): We present a neural analysis and synthesis (NANSY) framework that can manipulate the voice, pitch, and speed of an arbitrary speech signal. Most of the previous works have focused on using information bottleneck to disentangle analysis features for controllable synthesis, which usually results in poor reconstruction quality. We address this issue by proposing a novel training strategy based on information perturbation. The idea is to perturb information in the original input signal (e.g., formant, pitch, and frequency response), thereby letting synthesis networks selectively take essential attributes to reconstruct the input signal. Because NANSY does not need any bottleneck structures, it enjoys both high reconstruction quality and controllability. Furthermore, NANSY does not require any labels associated with speech data such as text and speaker information, but rather uses a new set of analysis features, i.e., wav2vec feature and newly proposed pitch feature, Yingram, which allows for fully self-supervised training. Taking advantage of fully self-supervised training, NANSY can be easily extended to a multilingual setting by simply training it with a multilingual dataset. The experiments show that NANSY can achieve significant improvement in performance in several applications such as zero-shot voice conversion, pitch shift, and time-scale modification.

[Auto-Encoding Knowledge Graph for Unsupervised Medical Report Generation](#)

- Fenglin Liu, Chenyu You, Xian Wu, Shen Ge, Sheng wang, Xu Sun
- abstract@[open-review](#): Medical report generation, which aims to automatically generate a long and coherent report of a given medical image, has been receiving growing research interests. Existing approaches mainly adopt a supervised manner and heavily rely on coupled image-report pairs. However, in the medical domain, building a large-scale image-report paired dataset is both time-consuming and expensive. To relax the dependency on paired data, we propose an unsupervised model Knowledge Graph Auto-Encoder (KGAE) which accepts independent sets of images and reports in training. KGAE consists of a pre-constructed knowledge graph, a knowledge-driven encoder and a knowledge-driven decoder. The knowledge graph works as the shared latent space to bridge the visual and textual domains; The knowledge-driven encoder projects medical images and reports to the corresponding coordinates in this latent space and the knowledge-driven decoder generates a medical report given a coordinate in this space. Since the knowledge-driven encoder and decoder can be trained with independent sets of images and reports, KGAE is unsupervised. The experiments show that the unsupervised KGAE generates desirable medical reports without using any image-report training pairs. Moreover, KGAE can also work in both semi-supervised and supervised settings, and accept paired images and reports in training. By further fine-tuning with image-report pairs, KGAE consistently outperforms the current state-of-the-art models on two datasets.

[Diffusion Normalizing Flow](#)

- Qinsheng Zhang, Yongxin Chen
- abstract@[open-review](#): We present a novel generative modeling method called diffusion normalizing flow based on stochastic differential equations (SDEs). The algorithm consists of two neural SDEs: a forward SDE that gradually adds noise to the data to transform the data into Gaussian random noise, and a backward SDE that gradually removes the noise to sample from the data distribution. By jointly training the two neural SDEs to minimize a common cost function that quantifies the difference between the two, the backward SDE converges to a diffusion process that starts with a Gaussian distribution and ends with the desired data distribution. Our method is closely related to normalizing flow and diffusion probabilistic models, and can be viewed as a combination of the two. Compared with normalizing flow, diffusion normalizing flow is able to learn distributions with sharp boundaries. Compared with diffusion probabilistic models, diffusion normalizing flow requires fewer discretization steps and thus has better sampling efficiency. Our algorithm demonstrates competitive performance in both high-dimension data density estimation and image generation tasks.

[Introspective Distillation for Robust Question Answering](#)

- Yulei Niu, Hanwang Zhang
- abstract@[open-review](#): Question answering (QA) models are well-known to exploit data bias, e.g., the language prior in visual QA and the position bias in reading comprehension. Recent debiasing methods achieve good out-of-distribution (OOD) generalizability with a considerable sacrifice of the in-distribution (ID) performance. Therefore, they are only applicable in domains where the test distribution is known in advance. In this paper, we present a novel debiasing method called Introspective Distillation (IntroD) to make the best of both worlds for QA. Our key technical contribution is to blend the inductive bias of OOD and ID by introspecting whether a training sample fits in the factual ID world or the counterfactual OOD one. Experiments on visual QA datasets VQA v2, VQA-CP, and reading comprehension dataset SQuAD demonstrate that our proposed IntroD maintains the competitive OOD performance compared to other debiasing methods, while sacrificing little or even achieving better ID performance compared to the non-debiasing ones.

[Rethinking the Pruning Criteria for Convolutional Neural Network](#)

- Zhongzhan Huang, Wenqi Shao, Xinjiang Wang, Liang Lin, Ping Luo
- abstract@[open-review](#): Channel pruning is a popular technique for compressing convolutional neural networks (CNNs), where various pruning criteria have been proposed to remove the redundant filters. From our comprehensive experiments, we found two blind spots of pruning criteria: (1) Similarity: There are some strong similarities among several primary pruning criteria that are widely cited and compared. According to these criteria, the ranks of filters' Importance Score are almost identical, resulting in similar pruned structures. (2) Applicability: The filters' Importance Score measured by some pruning criteria are too close to distinguish the network redundancy well. In this paper, we analyze the above blind spots on different types of pruning criteria with layer-wise pruning or global pruning. We also break some stereotypes, such as that the results of \$ell_1\$ and \$ell_2\$ pruning are not always similar. These analyses are based on the empirical experiments and our assumption (Convolutional Weight Distribution Assumption) that the well-trained convolutional filters in each layer approximately follow a Gaussian-alike distribution. This assumption has been verified through systematic and extensive statistical tests.

[Adaptive Machine Unlearning](#)

- Varun Gupta, Christopher Jung, Seth Neel, Aaron Roth, Saeed Sharifi-Malvajerdi, Chris Waites
- abstract@[open-review](#): Data deletion algorithms aim to remove the influence of deleted data points from trained models at a cheaper computational cost than fully retraining those models. However, for sequences of deletions, most prior work in the non-convex setting gives valid guarantees only for sequences that are chosen independently of the models that are published. If people choose to delete their data as a function of the published models (because they don't like what the models reveal about them, for example), then the update sequence is adaptive. In this paper, we give a general reduction from deletion guarantees against adaptive sequences to deletion guarantees against non-adaptive sequences, using differential privacy and its connection to max information. Combined with ideas from prior work which give guarantees for non-adaptive deletion sequences, this leads to extremely flexible algorithms able to handle arbitrary model classes and training methodologies, giving strong provable deletion guarantees for adaptive deletion sequences. We show in theory how prior work for non-convex models fails against adaptive deletion sequences, and use this intuition to design a practical attack against the SISA algorithm of Bourtoule et al. [2021] on CIFAR-10, MNIST, Fashion-MNIST.

[EditGAN: High-Precision Semantic Image Editing](#)

- Huan Ling, Karsten Kreis, Daiqing Li, Seung Wook Kim, Antonio Torralba, Sanja Fidler
- abstract@[open-review](#): Generative adversarial networks (GANs) have recently found applications in image editing. However, most GAN-based image editing methods often require large-scale datasets with semantic segmentation annotations for training, only provide high-level control, or merely interpolate between different images. Here, we propose EditGAN, a novel method for high-quality, high-precision semantic image editing, allowing users to edit images by modifying their highly detailed part segmentation masks, e.g., drawing a new mask for the headlight of a car. EditGAN builds on a GAN framework that jointly models images and their semantic segmentation, requiring only a handful of labeled examples – making it a scalable tool for editing. Specifically, we embed an image into the GAN's latent space and perform conditional latent code optimization according to the segmentation edit, which effectively also modifies the image. To amortize optimization, we find “editing vectors” in latent space that realize the edits. The framework allows us to learn an arbitrary number of editing vectors, which can then be directly applied on other images at interactive rates. We experimentally show that EditGAN can manipulate images with an unprecedented level of detail and freedom while preserving full image quality. We can also easily combine multiple edits and perform plausible edits beyond EditGAN's training data. We demonstrate EditGAN on a wide variety of image types and quantitatively outperform several previous editing methods on standard editing benchmark tasks.

[Deep Molecular Representation Learning via Fusing Physical and Chemical Information](#)

- Shuwen Yang, Ziyao Li, Guojie Song, Lingsheng Cai
- abstract@[open-review](#): Molecular representation learning is the first yet vital step in combining deep learning and molecular science. To push the boundaries of molecular representation learning, we present PhysChem, a novel neural architecture that learns molecular representations via fusing physical and chemical information of molecules. PhysChem is composed of a physicist network (PhysNet) and a chemist network (ChemNet). PhysNet is a neural physical engine that learns molecular conformations through simulating molecular dynamics with parameterized forces; ChemNet implements geometry-aware deep message-passing to learn chemical / biomedical properties of molecules. Two networks specialize in their own tasks and cooperate by providing expertise to each other. By fusing physical and chemical information, PhysChem achieved state-of-the-art performances on MoleculeNet, a standard molecular machine learning benchmark. The effectiveness of PhysChem was further corroborated on cutting-edge datasets of SARS-CoV-2.

[Neural optimal feedback control with local learning rules](#)

- Johannes Friedrich, Siavash Golkar, Shiva Farashahi, Alexander Genkin, Anirvan Sengupta, Dmitri Chklovskii
- abstract@[open-review](#): A major problem in motor control is understanding how the brain plans and executes proper movements in the face of delayed and noisy stimuli. A prominent framework for addressing such control problems is Optimal Feedback Control (OFC). OFC generates control actions that optimize behaviorally relevant criteria by integrating noisy sensory stimuli and the predictions of an internal model using the Kalman filter or its extensions. However, a satisfactory neural model of Kalman filtering and control is lacking because existing proposals have the following limitations: not considering the delay of sensory feedback, training in alternating phases, requiring knowledge of the noise covariance matrices, as well as that of systems dynamics. Moreover, the majority of these studies considered Kalman filtering in isolation, and not jointly with control. To address these shortcomings, we introduce a novel online algorithm which combines adaptive Kalman filtering with a model free control approach (i.e., policy gradient algorithm). We implement this algorithm in a biologically plausible neural network with local synaptic plasticity rules. This network, with local synaptic plasticity rules, performs system identification, Kalman filtering and control with delayed noisy sensory feedback. This network performs system identification and Kalman filtering, without the need for multiple phases with distinct update rules or the knowledge of the noise covariances. It can perform state estimation with delayed sensory feedback, with the help of an internal model. It learns the control policy without requiring any knowledge of the dynamics, thus avoiding the need for weight transport. In this way, our implementation of OFC solves the credit assignment problem needed to produce the appropriate sensory-motor control in the presence of stimulus delay.

[Reinforcement Learning in Linear MDPs: Constant Regret and Representation Selection](#)

- Matteo Papini, Andrea Tirinzoni, Aldo Pacchiano, Marcello Restelli, Alessandro Lazaric, Matteo Pirotta

- abstract@[open-review](#): We study the role of the representation of state-action value functions in regret minimization in finite-horizon Markov Decision Processes (MDPs) with linear structure. We first derive a necessary condition on the representation, called universally spanning optimal features (UNISOFT), to achieve constant regret in any MDP with linear reward function. This result encompasses the well-known settings of low-rank MDPs and, more generally, zero inherent Bellman error (also known as the Bellman closure assumption). We then demonstrate that this condition is also sufficient for these classes of problems by deriving a constant regret bound for two optimistic algorithms (LSVI-UCB and ELEANOR). Finally, we propose an algorithm for representation selection and we prove that it achieves constant regret when one of the given representations, or a suitable combination of them, satisfies the UNISOFT condition.

Noether Networks: meta-learning useful conserved quantities

- Ferran Alet, Dylan Doblar, Allan Zhou, Josh Tenenbaum, Kenji Kawaguchi, Chelsea Finn
- abstract@[open-review](#): Progress in machine learning (ML) stems from a combination of data availability, computational resources, and an appropriate encoding of inductive biases. Useful biases often exploit symmetries in the prediction problem, such as convolutional networks relying on translation equivariance. Automatically discovering these useful symmetries holds the potential to greatly improve the performance of ML systems, but still remains a challenge. In this work, we focus on sequential prediction problems and take inspiration from Noether's theorem to reduce the problem of finding inductive biases to meta-learning useful conserved quantities. We propose Noether Networks: a new type of architecture where a meta-learned conservation loss is optimized inside the prediction function. We show, theoretically and experimentally, that Noether Networks improve prediction quality, providing a general framework for discovering inductive biases in sequential problems.

Uncertainty-Driven Loss for Single Image Super-Resolution

- Qian Ning, Weisheng Dong, Xin Li, Jinjian Wu, GUANGMING Shi
- abstract@[open-review](#): In low-level vision such as single image super-resolution (SISR), traditional MSE or L1 loss function treats every pixel equally with the assumption that the importance of all pixels is the same. However, it has been long recognized that texture and edge areas carry more important visual information than smooth areas in photographic images. How to achieve such spatial adaptation in a principled manner has been an open problem in both traditional model-based and modern learning-based approaches toward SISR. In this paper, we propose a new adaptive weighted loss for SISR to train deep networks focusing on challenging situations such as textured and edge pixels with high uncertainty. Specifically, we introduce variance estimation characterizing the uncertainty on a pixel-by-pixel basis into SISR solutions so the targeted pixels in a high-resolution image (mean) and their corresponding uncertainty (variance) can be learned simultaneously. Moreover, uncertainty estimation allows us to leverage conventional wisdom such as sparsity prior for regularizing SISR solutions. Ultimately, pixels with large certainty (e.g., texture and edge pixels) will be prioritized for SISR according to their importance to visual quality. For the first time, we demonstrate that such uncertainty-driven loss can achieve better results than MSE or L1 loss for a wide range of network architectures. Experimental results on three popular SISR networks show that our proposed uncertainty-driven loss has achieved better PSNR performance than traditional loss functions without any increased computation during testing. The code is available at <https://see.xidian.edu.cn/faculty/wsdong/Projects/UDL-SR.htm>

GradInit: Learning to Initialize Neural Networks for Stable and Efficient Training

- Chen Zhu, Renkun Ni, Zheng Xu, Kezhi Kong, W. Ronny Huang, Tom Goldstein
- abstract@[open-review](#): Innovations in neural architectures have fostered significant breakthroughs in language modeling and computer vision. Unfortunately, novel architectures often result in challenging hyper-parameter choices and training instability if the network parameters are not properly initialized. A number of architecture-specific initialization schemes have been proposed, but these schemes are not always portable to new architectures. This paper presents GradInit, an automated and architecture agnostic method for initializing neural networks. GradInit is based on a simple heuristic; the norm of each network layer is adjusted so that a single step of SGD or Adam with prescribed hyperparameters results in the smallest possible loss value. This adjustment is done by introducing a scalar multiplier variable in front of each parameter block, and then optimizing these variables using a simple numerical scheme. GradInit accelerates the convergence and test performance of many convolutional architectures, both with or without skip connections, and even without normalization layers. It also improves the stability of the original Transformer architecture for machine translation, enabling training it without learning rate warmup using either Adam or SGD under a wide range of learning rates and momentum coefficients. Code is available at <https://github.com/zhuchen03/gradinit>.

Capacity and Bias of Learned Geometric Embeddings for Directed Graphs

- Michael Boratko, Dongxu Zhang, Nicholas Monath, Luke Vilnis, Kenneth L Clarkson, Andrew McCallum
- abstract@[open-review](#): A wide variety of machine learning tasks such as knowledge base completion, ontology alignment, and multi-label classification can benefit from incorporating into learning differentiable representations of graphs or taxonomies. While vectors in Euclidean space can theoretically represent any graph, much recent work shows that alternatives such as complex, hyperbolic, order, or box embeddings have geometric properties better suited to modeling real-world graphs. Experimentally these gains are seen only in lower dimensions, however, with performance benefits diminishing in higher dimensions. In this work, we introduce a novel variant of box embeddings that uses a learned smoothing parameter to achieve better representational capacity than vector models in low dimensions, while also avoiding performance saturation common to other geometric models in high dimensions. Further, we present theoretical results that prove box embeddings can represent any DAG. We perform rigorous empirical evaluations of vector, hyperbolic, and region-based geometric representations on several families of synthetic and real-world directed graphs. Analysis of these results exposes correlations between different families of graphs, graph characteristics, model size, and embedding geometry, providing useful insights into the inductive biases of various differentiable graph representations.

Online Learning Of Neural Computations From Sparse Temporal Feedback

- Lukas Braun, Tim Vogels
- abstract@[open-review](#): Neuronal computations depend on synaptic connectivity and intrinsic electrophysiological properties. Synaptic connectivity determines which inputs from presynaptic neurons are integrated, while cellular properties determine how inputs are filtered over time. Unlike their biological counterparts, most computational approaches to learning in simulated neural networks are limited to changes in synaptic connectivity. However, if intrinsic parameters change, neural computations are altered drastically. Here, we include the parameters that determine the intrinsic properties, e.g., time constants and reset potential, into the learning paradigm. Using sparse feedback signals that indicate target spike times, and gradient-based parameter updates, we show that the intrinsic parameters can be learned along with the synaptic weights to produce specific input-output functions. Specifically, we use a teacher-student paradigm in which a randomly initialised leaky integrate-and-fire or resonate-and-fire neuron must recover the parameters of a teacher neuron. We show that complex temporal functions can be learned online and without backpropagation through time, relying on event-based updates only. Our results are a step towards online learning of neural computations from ungraded and unsigned sparse feedback signals with a biologically inspired learning mechanism.

Self-Supervised Learning with Data Augmentations Provably Isolates Content from Style

- Julius von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel Besserve, Francesco Locatello
- abstract@[open-review](#): Self-supervised representation learning has shown remarkable success in a number of domains. A common practice is to perform data augmentation via hand-crafted transformations intended to leave the semantics of the data invariant. We seek to understand the empirical success of

this approach from a theoretical perspective. We formulate the augmentation process as a latent variable model by postulating a partition of the latent representation into a content component, which is assumed invariant to augmentation, and a style component, which is allowed to change. Unlike prior work on disentanglement and independent component analysis, we allow for both nontrivial statistical and causal dependencies in the latent space. We study the identifiability of the latent representation based on pairs of views of the observations and prove sufficient conditions that allow us to identify the invariant content partition up to an invertible mapping in both generative and discriminative settings. We find numerical simulations with dependent latent variables are consistent with our theory. Lastly, we introduce Causal3DIdent, a dataset of high-dimensional, visually complex images with rich causal dependencies, which we use to study the effect of data augmentations performed in practice.

[Instance-Conditional Knowledge Distillation for Object Detection](#)

- Zijian Kang, Peizhen Zhang, Xiangyu Zhang, Jian Sun, Nanning Zheng
- abstract@[open-review](#): Knowledge distillation has shown great success in classification, however, it is still challenging for detection. In a typical image for detection, representations from different locations may have different contributions to detection targets, making the distillation hard to balance. In this paper, we propose a conditional distillation framework to distill the desired knowledge, namely knowledge that is beneficial in terms of both classification and localization for every instance. The framework introduces a learnable conditional decoding module, which retrieves information given each target instance as query. Specifically, we encode the condition information as query and use the teacher's representations as key. The attention between query and key is used to measure the contribution of different features, guided by a localization-recognition-sensitive auxiliary task. Extensive experiments demonstrate the efficacy of our method: we observe impressive improvements under various settings. Notably, we boost RetinaNet with ResNet-50 backbone from \$37.4\$ to \$40.7\$ mAP (\$+3.3\$) under \$1\times\$ schedule, that even surpasses the teacher (\$40.4\$ mAP) with ResNet-101 backbone under \$3\times\$ schedule. Code has been released on <https://github.com/megvii-research/ICD>.

[Self-Supervised Representation Learning on Neural Network Weights for Model Characteristic Prediction](#)

- Konstantin Schürholt, Dimche Kostadinov, Damian Borth
- abstract@[open-review](#): Self-Supervised Learning (SSL) has been shown to learn useful and information-preserving representations. Neural Networks (NNs) are widely applied, yet their weight space is still not fully understood. Therefore, we propose to use SSL to learn hyper-representations of the weights of populations of NNs. To that end, we introduce domain specific data augmentations and an adapted attention architecture. Our empirical evaluation demonstrates that self-supervised representation learning in this domain is able to recover diverse NN model characteristics. Further, we show that the proposed learned representations outperform prior work for predicting hyper-parameters, test accuracy, and generalization gap as well as transfer to out-of-distribution settings.

[Multimodal Virtual Point 3D Detection](#)

- Tianwei Yin, Xingyi Zhou, Philipp Krähenbühl
- abstract@[open-review](#): Lidar-based sensing drives current autonomous vehicles. Despite rapid progress, current Lidar sensors still lag two decades behind traditional color cameras in terms of resolution and cost. For autonomous driving, this means that large objects close to the sensors are easily visible, but far-away or small objects comprise only one measurement or two. This is an issue, especially when these objects turn out to be driving hazards. On the other hand, these same objects are clearly visible in onboard RGB sensors. In this work, we present an approach to seamlessly fuse RGB sensors into Lidar-based 3D recognition. Our approach takes a set of 2D detections to generate dense 3D virtual points to augment an otherwise sparse 3D point cloud. These virtual points naturally integrate into any standard Lidar-based 3D detectors along with regular Lidar measurements. The resulting multi-modal detector is simple and effective. Experimental results on the large-scale nuScenes dataset show that our framework improves a strong CenterPoint baseline by a significant \$6.6\$ mAP, and outperforms competing fusion approaches. Code and more visualizations are available at <https://tianweiy.github.io/mvp/>

[On Joint Learning for Solving Placement and Routing in Chip Design](#)

- Ruoyu Cheng, Junchi Yan
- abstract@[open-review](#): For its advantage in GPU acceleration and less dependency on human experts, machine learning has been an emerging tool for solving the placement and routing problems, as two critical steps in modern chip design flow. Being still in its early stage, there are several fundamental issues unresolved: scalability, reward design, and end-to-end learning paradigm etc. To achieve end-to-end placement learning, we first propose a joint learning method for the placement of macros and standard cells, by the integration of reinforcement learning with a gradient based optimization scheme. To further bridge the placement with the subsequent routing task, we also develop a joint learning approach via reinforcement learning. One key design in our (reinforcement) learning paradigm involves a multi-view embedding model to encode both global graph level and local node level information of the input macros. Moreover, the random network distillation is devised to encourage exploration. Experiments on public chip design benchmarks show that our method can effectively learn from experience and also provide high-quality intermediate placement for the post standard cell placement, within few hours for training.

[Learning with Algorithmic Supervision via Continuous Relaxations](#)

- Felix Petersen, Christian Borgelt, Hilde Kuehne, Oliver Deussen
- abstract@[open-review](#): The integration of algorithmic components into neural architectures has gained increased attention recently, as it allows training neural networks with new forms of supervision such as ordering constraints or silhouettes instead of using ground truth labels. Many approaches in the field focus on the continuous relaxation of a specific task and show promising results in this context. But the focus on single tasks also limits the applicability of the proposed concepts to a narrow range of applications. In this work, we build on those ideas to propose an approach that allows to integrate algorithms into end-to-end trainable neural network architectures based on a general approximation of discrete conditions. To this end, we relax these conditions in control structures such as conditional statements, loops, and indexing, so that resulting algorithms are smoothly differentiable. To obtain meaningful gradients, each relevant variable is perturbed via logistic distributions and the expectation value under this perturbation is approximated. We evaluate the proposed continuous relaxation model on four challenging tasks and show that it can keep up with relaxations specifically designed for each individual task.

[Differentiable Multiple Shooting Layers](#)

- Stefano Massaroli, Michael Poli, Sho Sonoda, Taiji Suzuki, Jinkyoo Park, Atsushi Yamashita, Hajime Asama
- abstract@[open-review](#): We detail a novel class of implicit neural models. Leveraging time-parallel methods for differential equations, Multiple Shooting Layers (MSLs) seek solutions of initial value problems via parallelizable root-finding algorithms. MSLs broadly serve as drop-in replacements for neural ordinary differential equations (Neural ODEs) with improved efficiency in number of function evaluations (NFEs) and wall-clock inference time. We develop the algorithmic framework of MSLs, analyzing the different choices of solution methods from a theoretical and computational perspective. MSLs are showcased in long horizon optimal control of ODEs and PDEs and as latent models for sequence generation. Finally, we investigate the speedups obtained through application of MSL inference in neural controlled differential equations (Neural CDEs) for time series classification of medical data.

[Global-aware Beam Search for Neural Abstractive Summarization](#)

- Ye Ma, Zixun Lan, Lu Zong, Kaizhu Huang
- abstract@[open-review](#): This study develops a calibrated beam-based algorithm with awareness of the global attention distribution for neural abstractive summarization, aiming to improve the local optimality problem of the original beam search in a rigorous way. Specifically, a novel global protocol is proposed based on the attention distribution to stipulate how a global optimal hypothesis should attend to the source. A global scoring mechanism is then developed to regulate beam search to generate summaries in a near-global optimal fashion. This novel design enjoys a distinctive property, i.e., the global attention distribution could be predicted before inference, enabling step-wise improvements on the beam search through the global scoring mechanism. Extensive experiments on nine datasets show that the global (attention)-aware inference significantly improves state-of-the-art summarization models even using empirical hyper-parameters. The algorithm is also proven robust as it remains to generate meaningful texts with corrupted attention distributions. The codes and a comprehensive set of examples are available.

[DROID-SLAM: Deep Visual SLAM for Monocular, Stereo, and RGB-D Cameras](#)

- Zachary Teed, Jia Deng
- abstract@[open-review](#): We introduce DROID-SLAM, a new deep learning based SLAM system. DROID-SLAM consists of recurrent iterative updates of camera pose and pixelwise depth through a Dense Bundle Adjustment layer. DROID-SLAM is accurate, achieving large improvements over prior work, and robust, suffering from substantially fewer catastrophic failures. Despite training on monocular video, it can leverage stereo or RGB-D video to achieve improved performance at test time. The URL to our open source code is <https://github.com/princeton-vl/DROID-SLAM>.

[Few-Shot Object Detection via Association and DIscrimination](#)

- Yuhang Cao, Jiaqi Wang, Ying Jin, Tong Wu, Kai Chen, Ziwei Liu, Dahua Lin
- abstract@[open-review](#): Object detection has achieved substantial progress in the last decade. However, detecting novel classes with only few samples remains challenging, since deep learning under low data regime usually leads to a degraded feature space. Existing works employ a holistic fine-tuning paradigm to tackle this problem, where the model is first pre-trained on all base classes with abundant samples, and then it is used to carve the novel class feature space. Nonetheless, this paradigm is still imperfect. During fine-tuning, a novel class may implicitly leverage the knowledge of multiple base classes to construct its feature space, which induces a scattered feature space, hence violating the inter-class separability. To overcome these obstacles, we propose a two-step fine-tuning framework, Few-shot object detection via Association and DIscrimination (FADI), which builds up a discriminative feature space for each novel class with two integral steps. 1) In the association step, in contrast to implicitly leveraging multiple base classes, we construct a compact novel class feature space via explicitly imitating a specific base class feature space. Specifically, we associate each novel class with a base class according to their semantic similarity. After that, the feature space of a novel class can readily imitate the well-trained feature space of the associated base class. 2) In the discrimination step, to ensure the separability between the novel classes and associated base classes, we disentangle the classification branches for base and novel classes. To further enlarge the inter-class separability between all classes, a set-specialized margin loss is imposed. Extensive experiments on standard Pascal VOC and MS-COCO datasets demonstrate that FADI achieves new state-of-the-art performance, significantly improving the baseline in any shot/split by +18.7. Notably, the advantage of FADI is most announced on extremely few-shot scenarios (e.g. 1- and 3- shot).

[Neural Dubber: Dubbing for Videos According to Scripts](#)

- Chenxu Hu, Qiao Tian, Tingle Li, Wang Yuping, Yuxuan Wang, Hang Zhao
- abstract@[open-review](#): Dubbing is a post-production process of re-recording actors' dialogues, which is extensively used in filmmaking and video production. It is usually performed manually by professional voice actors who read lines with proper prosody, and in synchronization with the pre-recorded videos. In this work, we propose Neural Dubber, the first neural network model to solve a novel automatic video dubbing (AVD) task: synthesizing human speech synchronized with the given video from the text. Neural Dubber is a multi-modal text-to-speech (TTS) model that utilizes the lip movement in the video to control the prosody of the generated speech. Furthermore, an image-based speaker embedding (ISE) module is developed for the multi-speaker setting, which enables Neural Dubber to generate speech with a reasonable timbre according to the speaker's face. Experiments on the chemistry lecture single-speaker dataset and LRS2 multi-speaker dataset show that Neural Dubber can generate speech audios on par with state-of-the-art TTS models in terms of speech quality. Most importantly, both qualitative and quantitative evaluations show that Neural Dubber can control the prosody of synthesized speech by the video, and generate high-fidelity speech temporally synchronized with the video.

[Neural Bootstrapper](#)

- Minsuk Shin, Hyungjoo Cho, Hyun-seok Min, Sungbin Lim
- abstract@[open-review](#): Bootstrapping has been a primary tool for ensemble and uncertainty quantification in machine learning and statistics. However, due to its nature of multiple training and resampling, bootstrapping deep neural networks is computationally burdensome; hence it has difficulties in practical application to the uncertainty estimation and related tasks. To overcome this computational bottleneck, we propose a novel approach called Neural Bootstrapper (NeuBoots), which learns to generate bootstrapped neural networks through single model training. NeuBoots injects the bootstrap weights into the high-level feature layers of the backbone network and outputs the bootstrapped predictions of the target, without additional parameters and the repetitive computations from scratch. We apply NeuBoots to various machine learning tasks related to uncertainty quantification, including prediction calibrations in image classification and semantic segmentation, active learning, and detection of out-of-distribution samples. Our empirical results show that NeuBoots outperforms other bagging based methods under a much lower computational cost without losing the validity of bootstrapping.

[An Axiomatic Theory of Provably-Fair Welfare-Centric Machine Learning](#)

- Cyrus Cousins
- abstract@[open-review](#): We address an inherent difficulty in welfare-theoretic fair machine learning (ML), by proposing an equivalently-axiomatically justified alternative setting, and studying the resulting computational and statistical learning questions. Welfare metrics quantify overall wellbeing across a population of groups, and welfare-based objectives and constraints have recently been proposed to incentivize fair ML methods to satisfy their diverse needs. However, many ML problems are cast as loss minimization tasks, rather than utility maximization, and thus require nontrivial modeling to construct utility functions. We define a complementary metric, termed malfare, measuring overall societal harm, with axiomatic justification via the standard axioms of cardinal welfare, and cast fair ML as malfare minimization over the risk values (expected losses) of each group. Surprisingly, the axioms of cardinal welfare (malfare) dictate that this is not equivalent to simply defining utility as negative loss and maximizing welfare. Building upon these concepts, we define fair-PAC learning, where a fair-PAC learner is an algorithm that learns an ϵ - δ malfare-optimal model with bounded sample complexity, for any data distribution and (axiomatically justified) malfare concept. Finally, we show conditions under which many standard PAC-learners may be converted to fair-PAC learners, which places fair-PAC learning on firm theoretical ground, as it yields statistical — and in some cases computational — efficiency guarantees for many well-studied ML models. Fair-PAC learning is also practically relevant, as it democratizes fair ML by providing concrete training algorithms with rigorous generalization guarantees.

[HSVA: Hierarchical Semantic-Visual Adaptation for Zero-Shot Learning](#)

- Shiming Chen, Guosen Xie, Yang Liu, Qinmu Peng, Baigui Sun, Hao Li, Xinge You, Ling Shao
- abstract@[open-review](#): Zero-shot learning (ZSL) tackles the unseen class recognition problem, transferring semantic knowledge from seen classes to unseen ones. Typically, to guarantee desirable knowledge transfer, a common (latent) space is adopted for associating the visual and semantic domains in ZSL. However, existing common space learning methods align the semantic and visual domains by merely mitigating distribution disagreement through

one-step adaptation. This strategy is usually ineffective due to the heterogeneous nature of the feature representations in the two domains, which intrinsically contain both distribution and structure variations. To address this and advance ZSL, we propose a novel hierarchical semantic-visual adaptation (HSVA) framework. Specifically, HSVA aligns the semantic and visual domains by adopting a hierarchical two-step adaptation, i.e., structure adaptation and distribution adaptation. In the structure adaptation step, we take two task-specific encoders to encode the source data (visual domain) and the target data (semantic domain) into a structure-aligned common space. To this end, a supervised adversarial discrepancy (SAD) module is proposed to adversarially minimize the discrepancy between the predictions of two task-specific classifiers, thus making the visual and semantic feature manifolds more closely aligned. In the distribution adaptation step, we directly minimize the Wasserstein distance between the latent multivariate Gaussian distributions to align the visual and semantic distributions using a common encoder. Finally, the structure and distribution adaptation are derived in a unified framework under two partially-aligned variational autoencoders. Extensive experiments on four benchmark datasets demonstrate that HSVA achieves superior performance on both conventional and generalized ZSL. The code is available at \url{https://github.com/shiming-chen/HSVA}.

[Higher Order Kernel Mean Embeddings to Capture Filtrations of Stochastic Processes](#)

- Christopher Salvi, Maud Lemercier, Chong Liu, Blanka Horvath, Theodoros Damoulas, Terry Lyons
- abstract@[open-review](#): Stochastic processes are random variables with values in some space of paths. However, reducing a stochastic process to a path-valued random variable ignores its filtration, i.e. the flow of information carried by the process through time. By conditioning the process on its filtration, we introduce a family of higher order kernel mean embeddings (KMEs) that generalizes the notion of KME to capture additional information related to the filtration. We derive empirical estimators for the associated higher order maximum mean discrepancies (MMDs) and prove consistency. We then construct a filtration-sensitive kernel two-sample test able to capture information that gets missed by the standard MMD test. In addition, leveraging our higher order MMDs we construct a family of universal kernels on stochastic processes that allows to solve real-world calibration and optimal stopping problems in quantitative finance (such as the pricing of American options) via classical kernel-based regression methods. Finally, adapting existing tests for conditional independence to the case of stochastic processes, we design a causal-discovery algorithm to recover the causal graph of structural dependencies among interacting bodies solely from observations of their multidimensional trajectories.

[Low-Rank Subspaces in GANs](#)

- Jiapeng Zhu, Ruili Feng, Yujun Shen, Deli Zhao, Zheng-Jun Zha, Jingren Zhou, Qifeng Chen
- abstract@[open-review](#): The latent space of a Generative Adversarial Network (GAN) has been shown to encode rich semantics within some subspaces. To identify these subspaces, researchers typically analyze the statistical information from a collection of synthesized data, and the identified subspaces tend to control image attributes globally (i.e., manipulating an attribute causes the change of an entire image). By contrast, this work introduces low-rank subspaces that enable more precise control of GAN generation. Concretely, given an arbitrary image and a region of interest (e.g., eyes of face images), we manage to relate the latent space to the image region with the Jacobian matrix and then use low-rank factorization to discover steerable latent subspaces. There are three distinguishable strengths of our approach that can be aptly called LowRankGAN. First, compared to analytic algorithms in prior work, our low-rank factorization of Jacobians is able to find the low-dimensional representation of attribute manifold, making image editing more precise and controllable. Second, low-rank factorization naturally yields a null space of attributes such that moving the latent code within it only affects the outer region of interest. Therefore, local image editing can be simply achieved by projecting an attribute vector into the null space without relying on a spatial mask as existing methods do. Third, our method can robustly work with a local region from one image for analysis yet well generalize to other images, making it much easy to use in practice. Extensive experiments on state-of-the-art GAN models (including StyleGAN2 and BigGAN) trained on various datasets demonstrate the effectiveness of our LowRankGAN.

[Neural Symplectic Form: Learning Hamiltonian Equations on General Coordinate Systems](#)

- Yuhan Chen, Takashi Matsubara, Takaharu Yaguchi
- abstract@[open-review](#): In recent years, substantial research on the methods for learning Hamiltonian equations has been conducted. Although these approaches are very promising, the commonly used representation of the Hamilton equation uses the generalized momenta, which are generally unknown. Therefore, the training data must be represented in this unknown coordinate system, and this causes difficulty in applying the model to real data. Meanwhile, Hamiltonian equations also have a coordinate-free expression that is expressed by using the symplectic 2-form. In this study, we propose a model that learns the symplectic form from data using neural networks, thereby providing a method for learning Hamiltonian equations from data represented in general coordinate systems, which are not limited to the generalized coordinates and the generalized momenta. Consequently, the proposed method is capable not only of modeling target equations of both Hamiltonian and Lagrangian formalisms but also of extracting unknown Hamiltonian structures hidden in the data. For example, many polynomial ordinary differential equations such as the Lotka-Volterra equation are known to admit non-trivial Hamiltonian structures, and our numerical experiments show that such structures can be certainly learned from data. Technically, each symplectic 2-form is associated with a skew-symmetric matrix, but not all skew-symmetric matrices define the symplectic 2-form. In the proposed method, using the fact that symplectic 2-forms are derived as the exterior derivative of certain differential 1-forms, we model the differential 1-form by neural networks, thereby improving the efficiency of learning.

[Sample-Efficient Reinforcement Learning Is Feasible for Linearly Realizable MDPs with Limited Revisiting](#)

- Gen Li, Yuxin Chen, Yuejie Chi, Yuantao Gu, Yuting Wei
- abstract@[open-review](#): Low-complexity models such as linear function representation play a pivotal role in enabling sample-efficient reinforcement learning (RL). The current paper pertains to a scenario with value-based linear representation, which postulates linear realizability of the optimal Q-function (also called the ``linear Q^{\star} problem''). While linear realizability alone does not allow for sample-efficient solutions in general, the presence of a large sub-optimality gap is a potential game changer, depending on the sampling mechanism in use. Informally, sample efficiency is achievable with a large sub-optimality gap when a generative model is available, but is unfortunately infeasible when we turn to standard online RL settings. We make progress towards understanding this linear Q^{\star} problem by investigating a new sampling protocol, which draws samples in an online/exploratory fashion but allows one to backtrack and revisit previous states. This protocol is more flexible than the standard online RL setting, while being practically relevant and far more restrictive than the generative model. We develop an algorithm tailored to this setting, achieving a sample complexity that scales polynomially with the feature dimension, the horizon, and the inverse sub-optimality gap, but not the size of the state/action space. Our findings underscore the fundamental interplay between sampling protocols and low-complexity function representation in RL.

[Self-Paced Contrastive Learning for Semi-supervised Medical Image Segmentation with Meta-labels](#)

- Jizong Peng, Ping Wang, Christian Desrosiers, Marco Pedersoli
- abstract@[open-review](#): The contrastive pre-training of a recognition model on a large dataset of unlabeled data often boosts the model's performance on downstream tasks like image classification. However, in domains such as medical imaging, collecting unlabeled data can be challenging and expensive. In this work, we consider the task of medical image segmentation and adapt contrastive learning with meta-label annotations to scenarios where no additional unlabeled data is available. Meta-labels, such as the location of a 2D slice in a 3D MRI scan, often come for free during the acquisition process. We use these meta-labels to pre-train the image encoder, as well as in a semi-supervised learning step that leverages a reduced set of annotated data. A self-paced learning strategy exploiting the weak annotations is proposed to further help the learning process and discriminate useful labels from noise. Results on five medical image segmentation datasets show that our approach: i) highly boosts the performance of a model trained on a few scans, ii) outperforms previous contrastive and semi-supervised approaches, and iii) reaches close to the performance of a model trained on the full data.

[Reverse engineering recurrent neural networks with Jacobian switching linear dynamical systems](#)

- Jimmy Smith, Scott Linderman, David Sussillo
- abstract@[open-review](#): Recurrent neural networks (RNNs) are powerful models for processing time-series data, but it remains challenging to understand how they function. Improving this understanding is of substantial interest to both the machine learning and neuroscience communities. The framework of reverse engineering a trained RNN by linearizing around its fixed points has provided insight, but the approach has significant challenges. These include difficulty choosing which fixed point to expand around when studying RNN dynamics and error accumulation when reconstructing the nonlinear dynamics with the linearized dynamics. We present a new model that overcomes these limitations by co-training an RNN with a novel switching linear dynamical system (SLDS) formulation. A first-order Taylor series expansion of the co-trained RNN and an auxiliary function trained to pick out the RNN's fixed points govern the SLDS dynamics. The results are a trained SLDS variant that closely approximates the RNN, an auxiliary function that can produce a fixed point for each point in state-space, and a trained nonlinear RNN whose dynamics have been regularized such that its first-order terms perform the computation, if possible. This model removes the post-training fixed point optimization and allows us to unambiguously study the learned dynamics of the SLDS at any point in state-space. It also generalizes SLDS models to continuous manifolds of switching points while sharing parameters across switches. We validate the utility of the model on two synthetic tasks relevant to previous work reverse engineering RNNs. We then show that our model can be used as a drop-in in more complex architectures, such as LFADS, and apply this LFADS hybrid to analyze single-trial spiking activity from the motor system of a non-human primate.

[Learning-Augmented Dynamic Power Management with Multiple States via New Ski Rental Bounds](#)

- Antonios Antoniadis, Christian Coester, Marek Elias, Adam Polak, Bertrand Simon
- abstract@[open-review](#): We study the online problem of minimizing power consumption in systems with multiple power-saving states. During idle periods of unknown lengths, an algorithm has to choose between power-saving states of different energy consumption and wake-up costs. We develop a learning-augmented online algorithm that makes decisions based on (potentially inaccurate) predicted lengths of the idle periods. The algorithm's performance is near-optimal when predictions are accurate and degrades gracefully with increasing prediction error, with a worst-case guarantee almost identical to the optimal classical online algorithm for the problem. A key ingredient in our approach is a new algorithm for the online ski-rental problem in the learning augmented setting with tight dependence on the prediction error. We support our theoretical findings with experiments.

[Learning Equivariant Energy Based Models with Equivariant Stein Variational Gradient Descent](#)

- Priyank Jaini, Lars Holdijk, Max Welling
- abstract@[open-review](#): We focus on the problem of efficient sampling and learning of probability densities by incorporating symmetries in probabilistic models. We first introduce Equivariant Stein Variational Gradient Descent algorithm -- an equivariant sampling method based on Stein's identity for sampling from densities with symmetries. Equivariant SVGD explicitly incorporates symmetry information in a density through equivariant kernels which makes the resultant sampler efficient both in terms of sample complexity and the quality of generated samples. Subsequently, we define equivariant energy based models to model invariant densities that are learned using contrastive divergence. By utilizing our equivariant SVGD for training equivariant EBMs, we propose new ways of improving and scaling up training of energy based models. We apply these equivariant energy models for modelling joint densities in regression and classification tasks for image datasets, many-body particle systems and molecular structure generation.

[Information Directed Sampling for Sparse Linear Bandits](#)

- Botao Hao, Tor Lattimore, Wei Deng
- abstract@[open-review](#): Stochastic sparse linear bandits offer a practical model for high-dimensional online decision-making problems and have a rich information-regret structure. In this work we explore the use of information-directed sampling (IDS), which naturally balances the information-regret trade-off. We develop a class of information-theoretic Bayesian regret bounds that nearly match existing lower bounds on a variety of problem instances, demonstrating the adaptivity of IDS. To efficiently implement sparse IDS, we propose an empirical Bayesian approach for sparse posterior sampling using a spike-and-slab Gaussian-Laplace prior. Numerical results demonstrate significant regret reductions by sparse IDS relative to several baselines.

[Linear Convergence of Gradient Methods for Estimating Structured Transition Matrices in High-dimensional Vector Autoregressive Models](#)

- Xiao Lv, Wei Cui, Yulong Liu
- abstract@[open-review](#): In this paper, we present non-asymptotic optimization guarantees of gradient descent methods for estimating structured transition matrices in high-dimensional vector autoregressive (VAR) models. We adopt the projected gradient descent (PGD) for single-structured transition matrices and the alternating projected gradient descent (AltPGD) for superposition-structured ones. Our analysis demonstrates that both gradient algorithms converge linearly to the statistical error even though the strong convexity of the objective function is absent under the high-dimensional settings. Moreover our result is sharp (up to a constant factor) in the sense of matching the phase transition theory of the corresponding model with independent samples. To the best of our knowledge, this analysis constitutes first non-asymptotic optimization guarantees of the linear rate for regularized estimation in high-dimensional VAR models. Numerical results are provided to support our theoretical analysis.

[Large-Scale Unsupervised Object Discovery](#)

- Van Huy Vo, Elena Sizikova, Cordelia Schmid, Patrick Pérez, Jean Ponce
- abstract@[open-review](#): Existing approaches to unsupervised object discovery (UOD) do not scale up to large datasets without approximations that compromise their performance. We propose a novel formulation of UOD as a ranking problem, amenable to the arsenal of distributed methods available for eigenvalue problems and link analysis. Through the use of self-supervised features, we also demonstrate the first effective fully unsupervised pipeline for UOD. Extensive experiments on COCO~\cite{Lin2014cocodataset} and OpenImages~\cite{openimages} show that, in the single-object discovery setting where a single prominent object is sought in each image, the proposed LOD (Large-scale Object Discovery) approach is on par with, or better than the state of the art for medium-scale datasets (up to 120K images), and over 37% better than the only other algorithms capable of scaling up to 1.7M images. In the multi-object discovery setting where multiple objects are sought in each image, the proposed LOD is over 14% better in average precision (AP) than all other methods for datasets ranging from 20K to 1.7M images. Using self-supervised features, we also show that the proposed method obtains state-of-the-art UOD performance on OpenImages.

[Sparse Steerable Convolutions: An Efficient Learning of SE\(3\)-Equivariant Features for Estimation and Tracking of Object Poses in 3D Space](#)

- Jiehong Lin, Hongyang Li, Ke Chen, Jiangbo Lu, Kui Jia
- abstract@[open-review](#): As a basic component of SE(3)-equivariant deep feature learning, steerable convolution has recently demonstrated its advantages for 3D semantic analysis. The advantages are, however, brought by expensive computations on dense, volumetric data, which prevent its practical use for efficient processing of 3D data that are inherently sparse. In this paper, we propose a novel design of Sparse Steerable Convolution (SS-Conv) to address the shortcoming; SS-Conv greatly accelerates steerable convolution with sparse tensors, while strictly preserving the property of SE(3)-equivariance. Based on SS-Conv, we propose a general pipeline for precise estimation of object poses, wherein a key design is a Feature-Steering module that takes the

full advantage of SE(3)-equivariance and is able to conduct an efficient pose refinement. To verify our designs, we conduct thorough experiments on three tasks of 3D object semantic analysis, including instance-level 6D pose estimation, category-level 6D pose and size estimation, and category-level 6D pose tracking. Our proposed pipeline based on SS-Conv outperforms existing methods on almost all the metrics evaluated by the three tasks. Ablation studies also show the superiority of our SS-Conv over alternative convolutions in terms of both accuracy and efficiency. Our code is released publicly at <https://github.com/Gorilla-Lab-SCUT/SS-Conv>.

[Noisy Adaptation Generates Lévy Flights in Attractor Neural Networks](#)

- Xingsi Dong, Tianhao Chu, Tiejun Huang, Zilong Ji, Si Wu
- abstract@[open-review](#): Lévy flights describe a special class of random walks whose step sizes satisfy a power-law tailed distribution. As being an efficient searching strategy in unknown environments, Lévy flights are widely observed in animal foraging behaviors. Recent studies further showed that human cognitive functions also exhibit the characteristics of Lévy flights. Despite being a general phenomenon, the neural mechanism at the circuit level for generating Lévy flights remains unresolved. Here, we investigate how Lévy flights can be achieved in attractor neural networks. To elucidate the underlying mechanism clearly, we first study continuous attractor neural networks (CANNs), and find that noisy neural adaptation, exemplified by spike frequency adaptation (SFA) in this work, can generate Lévy flights representing transitions of the network state in the attractor space. Specifically, the strength of SFA defines a travelling wave boundary, below which the network state displays local Brownian motion, and above which the network state displays long-jump motion. Noises in neural adaptation causes the network state to intermittently switch between these two motion modes, manifesting the characteristics of Lévy flights. We further extend the study to a general attractor neural network, and demonstrate that our model can explain the Lévy-flight phenomenon observed during free memory retrieval of humans. We hope that this study will give us insight into understanding the neural mechanism for optimal information processing in the brain.

[On Linear Stability of SGD and Input-Smoothness of Neural Networks](#)

- Chao Ma, Lexing Ying
- abstract@[open-review](#): The multiplicative structure of parameters and input data in the first layer of neural networks is explored to build connection between the landscape of the loss function with respect to parameters and the landscape of the model function with respect to input data. By this connection, it is shown that flat minima regularize the gradient of the model function, which explains the good generalization performance of flat minima. Then, we go beyond the flatness and consider high-order moments of the gradient noise, and show that Stochastic Gradient Descent (SGD) tends to impose constraints on these moments by a linear stability analysis of SGD around global minima. Together with the multiplicative structure, we identify the Sobolev regularization effect of SGD, i.e. SGD regularizes the Sobolev seminorms of the model function with respect to the input data. Finally, bounds for generalization error and adversarial robustness are provided for solutions found by SGD under assumptions of the data distribution.

[Joint inference and input optimization in equilibrium networks](#)

- Swaminathan Gurumurthy, Shaojie Bai, Zachary Manchester, J. Zico Kolter
- abstract@[open-review](#): Many tasks in deep learning involve optimizing over the inputs to a network to minimize or maximize some objective; examples include optimization over latent spaces in a generative model to match a target image, or adversarially perturbing an input to worsen classifier performance. Performing such optimization, however, is traditionally quite costly, as it involves a complete forward and backward pass through the network for each gradient step. In a separate line of work, a recent thread of research has developed the deep equilibrium (DEQ) model, a class of models that foregoes traditional network depth and instead computes the output of a network by finding the fixed point of a single nonlinear layer. In this paper, we show that there is a natural synergy between these two settings. Although, naively using DEQs for these optimization problems is expensive (owing to the time needed to compute a fixed point for each gradient step), we can leverage the fact that gradient-based optimization can itself be cast as a fixed point iteration to substantially improve the overall speed. That is, we simultaneously both solve for the DEQ fixed point and optimize over network inputs, all within a single "augmented" DEQ model that jointly encodes both the original network and the optimization process. Indeed, the procedure is fast enough that it allows us to efficiently train DEQ models for tasks traditionally relying on an "inner" optimization loop. We demonstrate this strategy on various tasks such as training generative models while optimizing over latent codes, training models for inverse problems like denoising and inpainting, adversarial training and gradient based meta-learning.

[A unified framework for bandit multiple testing](#)

- Ziyu Xu, Ruodu Wang, Aaditya Ramdas
- abstract@[open-review](#): In bandit multiple hypothesis testing, each arm corresponds to a different null hypothesis that we wish to test, and the goal is to design adaptive algorithms that correctly identify large set of interesting arms (true discoveries), while only mistakenly identifying a few uninteresting ones (false discoveries). One common metric in non-bandit multiple testing is the false discovery rate (FDR). We propose a unified, modular framework for bandit FDR control that emphasizes the decoupling of exploration and summarization of evidence. We utilize the powerful martingale-based concept of "e-processes" to ensure FDR control for arbitrary composite nulls, exploration rules and stopping times in generic problem settings. In particular, valid FDR control holds even if the reward distributions of the arms could be dependent, multiple arms may be queried simultaneously, and multiple (cooperating or competing) agents may be querying arms, covering combinatorial semi-bandit type settings as well. Prior work has considered in great detail the setting where each arm's reward distribution is independent and sub-Gaussian, and a single arm is queried at each step. Our framework recovers matching sample complexity guarantees in this special case, and performs comparably or better in practice. For other settings, sample complexities will depend on the finer details of the problem (composite nulls being tested, exploration algorithm, data dependence structure, stopping rule) and we do not explore these; our contribution is to show that the FDR guarantee is clean and entirely agnostic to these details.

[Recovering Latent Causal Factor for Generalization to Distributional Shifts](#)

- Xinwei Sun, Botong Wu, Xiangyu Zheng, Chang Liu, Wei Chen, Tao Qin, Tie-Yan Liu
- abstract@[open-review](#): Distributional shifts between training and target domains may degrade the prediction accuracy of learned models, mainly because these models often learn features that possess only correlation rather than causal relation with the output. Such a correlation, which is known as ``spurious correlation'' statistically, is domain-dependent hence may fail to generalize to unseen domains. To avoid such a spurious correlation, we propose \textbf{La}tent \textbf{C}ausal \textbf{I}nvariance \textbf{M}odels (LaCIM) that specifies the underlying causal structure of the data and the source of distributional shifts, guiding us to pursue only causal factor for prediction. Specifically, the LaCIM introduces a pair of correlated latent factors: (a) causal factor and (b) others, while the extent of this correlation is governed by a domain variable that characterizes the distributional shifts. On the basis of this, we prove that the distribution of observed variables conditioning on latent variables is shift-invariant. Equipped with such an invariance, we prove that the causal factor can be recovered without mixing information from others, which induces the ground-truth predicting mechanism. We propose a Variational-Bayesian-based method to learn this invariance for prediction. The utility of our approach is verified by improved generalization to distributional shifts on various real-world data. Our code is freely available at \url{https://github.com/wubotong/LaCIM}.

[Graph Differentiable Architecture Search with Structure Learning](#)

- Yijian Qin, Xin Wang, Zeyang Zhang, Wenwu Zhu
- abstract@[open-review](#): Discovering ideal Graph Neural Networks (GNNs) architectures for different tasks is labor intensive and time consuming. To save human efforts, Neural Architecture Search (NAS) recently has been used to automatically discover adequate GNN architectures for certain tasks in order

to achieve competitive or even better performance compared with manually designed architectures. However, existing works utilizing NAS to search GNN structures fail to answer the question: how NAS is able to select the desired GNN architectures? In this paper, we investigate this question to solve the problem, for the first time. We conduct a measurement study with experiments to discover that gradient based NAS methods tend to select proper architectures based on the usefulness of different types of information with respect to the target task. Our explorations further show that gradient based NAS also suffers from noises hidden in the graph, resulting in searching suboptimal GNN architectures. Based on our findings, we propose a Graph differentiable Architecture Search model with Structure Optimization (GASSO), which allows differentiable search of the architecture with gradient descent and is able to discover graph neural architectures with better performance through employing graph structure learning as a denoising process in the search procedure. The proposed GASSO model is capable of simultaneously searching the optimal architecture and adaptively adjusting graph structure by jointly optimizing graph architecture search and graph structure denoising. Extensive experiments on real-world graph datasets demonstrate that our proposed GASSO model is able to achieve state-of-the-art performance compared with existing baselines.

[Designing Counterfactual Generators using Deep Model Inversion](#)

- Jayaraman Thiagarajan, Vivek Sivaraman Narayanaswamy, Deepta Rajan, Jia Liang, Akshay Chaudhari, Andreas Spanias
- abstract@[open-review](#): Explanation techniques that synthesize small, interpretable changes to a given image while producing desired changes in the model prediction have become popular for introspecting black-box models. Commonly referred to as counterfactuals, the synthesized explanations are required to contain discernible changes (for easy interpretability) while also being realistic (consistency to the data manifold). In this paper, we focus on the case where we have access only to the trained deep classifier and not the actual training data. While the problem of inverting deep models to synthesize images from the training distribution has been explored, our goal is to develop a deep inversion approach to generate counterfactual explanations for a given query image. Despite their effectiveness in conditional image synthesis, we show that existing deep inversion methods are insufficient for producing meaningful counterfactuals. We propose DISC (Deep Inversion for Synthesizing Counterfactuals) that improves upon deep inversion by utilizing (a) stronger image priors, (b) incorporating a novel manifold consistency objective and (c) adopting a progressive optimization strategy. We find that, in addition to producing visually meaningful explanations, the counterfactuals from DISC are effective at learning classifier decision boundaries and are robust to unknown test-time corruptions.

[A Faster Maximum Cardinality Matching Algorithm with Applications in Machine Learning](#)

- Nathaniel Lahn, Sharath Raghvendra, Jiacheng Ye
- abstract@[open-review](#): Maximum cardinality bipartite matching is an important graph optimization problem with several applications. For instance, maximum cardinality matching in a $\$delta$$ -disc graph can be used in the computation of the bottleneck matching as well as the $\$infty$$ -Wasserstein and the Lévy-Prokhorov distances between probability distributions. For any point sets $A, B \subset \mathbb{R}^2$, the $\$delta$$ -disc graph is a bipartite graph formed by connecting every pair of points $(a, b) \in A \times B$ by an edge if the Euclidean distance between them is at most $\$delta$$. Using the classical Hopcroft-Karp algorithm, a maximum-cardinality matching on any $\$delta$$ -disc graph can be found in $\tilde{O}(n^{3/2})$ time.^{~\footnote{We use $\tilde{O}(\cdot)$ to suppress poly-logarithmic terms in the complexity.}} In this paper, we present a simplification of a recent algorithm (Lahn and Raghvendra, JoCG 2021) for the maximum cardinality matching problem and describe how a maximum cardinality matching in a $\$delta$$ -disc graph can be computed asymptotically faster than $O(n^{3/2})$ time for any moderately dense point set. As applications, we show that if A and B are point sets drawn uniformly at random from a unit square, an exact bottleneck matching can be computed in $\tilde{O}(n^{4/3})$ time. On the other hand, experiments suggest that the Hopcroft-Karp algorithm seems to take roughly $\Theta(n^{3/2})$ time for this case. This translates to substantial improvements in execution time for larger inputs.

[Dynamic population-based meta-learning for multi-agent communication with natural language](#)

- Abhinav Gupta, Marc Lanctot, Angeliki Lazaridou
- abstract@[open-review](#): In this work, our goal is to train agents that can coordinate with seen, unseen as well as human partners in a multi-agent communication environment involving natural language. Previous work using a single set of agents has shown great progress in generalizing to known partners, however it struggles when coordinating with unfamiliar agents. To mitigate that, recent work explored the use of population-based approaches, where multiple agents interact with each other with the goal of learning more generic protocols. These methods, while able to result in good coordination between unseen partners, still only achieve so in cases of simple languages, thus failing to adapt to human partners using natural language. We attribute this to the use of static populations and instead propose a dynamic population-based meta-learning approach that builds such a population in an iterative manner. We perform a holistic evaluation of our method on two different referential games, and show that our agents outperform all prior work when communicating with seen partners and humans. Furthermore, we analyze the natural language generation skills of our agents, where we find that our agents also outperform strong baselines. Finally, we test the robustness of our agents when communicating with out-of-population agents and carefully test the importance of each component of our method through ablation studies.

[Adversarial Neuron Pruning Purifies Backdoored Deep Models](#)

- Dongxian Wu, Yisen Wang
- abstract@[open-review](#): As deep neural networks (DNNs) are growing larger, their requirements for computational resources become huge, which makes outsourcing training more popular. Training in a third-party platform, however, may introduce potential risks that a malicious trainer will return backdoored DNNs, which behave normally on clean samples but output targeted misclassifications whenever a trigger appears at the test time. Without any knowledge of the trigger, it is difficult to distinguish or recover benign DNNs from backdoored ones. In this paper, we first identify an unexpected sensitivity of backdoored DNNs, that is, they are much easier to collapse and tend to predict the target label on clean samples when their neurons are adversarially perturbed. Based on these observations, we propose a novel model repairing method, termed Adversarial Neuron Pruning (ANP), which prunes some sensitive neurons to purify the injected backdoor. Experiments show, even with only an extremely small amount of clean data (e.g., 1%), ANP effectively removes the injected backdoor without causing obvious performance degradation.

[Towards Robust and Reliable Algorithmic Recourse](#)

- Sohini Upadhyay, Shalmali Joshi, Himabindu Lakkaraju
- abstract@[open-review](#): As predictive models are increasingly being deployed in high-stakes decision making (e.g., loan approvals), there has been growing interest in post-hoc techniques which provide recourse to affected individuals. These techniques generate recourses under the assumption that the underlying predictive model does not change. However, in practice, models are often regularly updated for a variety of reasons (e.g., dataset shifts), thereby rendering previously prescribed recourses ineffective. To address this problem, we propose a novel framework, ROBust Algorithmic Recourse (ROAR), that leverages adversarial training for finding recourses that are robust to model shifts. To the best of our knowledge, this work proposes the first ever solution to this critical problem. We also carry out theoretical analysis which underscores the importance of constructing recourses that are robust to model shifts: 1) We quantify the probability of invalidation for recourses generated without accounting for model shifts. 2) We prove that the additional cost incurred due to the robust recourses output by our framework is bounded. Experimental evaluation on multiple synthetic and real-world datasets demonstrates the efficacy of the proposed framework.

[Neural Rule-Execution Tracking Machine For Transformer-Based Text Generation](#)

- Yufei Wang, Can Xu, Huang Hu, Chongyang Tao, Stephen Wan, Mark Dras, Mark Johnson, Dixin Jiang

- abstract@[open-review](#): Sequence-to-Sequence (Seq2Seq) neural text generation models, especially the pre-trained ones (e.g., BART and T5), have exhibited compelling performance on various natural language generation tasks. However, the black-box nature of these models limits their application in tasks where specific rules (e.g., controllable constraints, prior knowledge) need to be executed. Previous works either design specific model structures (e.g., Copy Mechanism corresponding to the rule "the generated output should include certain words in the source input") or implement specialized inference algorithms (e.g., Constrained Beam Search) to execute particular rules through the text generation. These methods require the careful design case-by-case and are difficult to support multiple rules concurrently. In this paper, we propose a novel module named Neural Rule-Execution Tracking Machine (NRETM) that can be equipped into various transformer-based generators to leverage multiple rules simultaneously to guide the neural generation model for superior generation performance in an unified and scalable way. Extensive experiments on several benchmarks verify the effectiveness of our proposed model in both controllable and general text generation tasks.

[Scalable Online Planning via Reinforcement Learning Fine-Tuning](#)

- Arnaud Fickinger, Hengyuan Hu, Brandon Amos, Stuart Russell, Noam Brown
- abstract@[open-review](#): Lookahead search has been a critical component of recent AI successes, such as in the games of chess, go, and poker. However, the search methods used in these games, and in many other settings, are tabular. Tabular search methods do not scale well with the size of the search space, and this problem is exacerbated by stochasticity and partial observability. In this work we replace tabular search with online model-based fine-tuning of a policy neural network via reinforcement learning, and show that this approach outperforms state-of-the-art search algorithms in benchmark settings. In particular, we use our search algorithm to achieve a new state-of-the-art result in self-play Hanabi, and show the generality of our algorithm by also showing that it outperforms tabular search in the Atari game Ms. Pacman.

[Adversarial Regression with Doubly Non-negative Weighting Matrices](#)

- Tam Le, Truyen Nguyen, Makoto Yamada, Jose Blanchet, Viet Anh Nguyen
- abstract@[open-review](#): Many machine learning tasks that involve predicting an output response can be solved by training a weighted regression model. Unfortunately, the predictive power of this type of models may severely deteriorate under low sample sizes or under covariate perturbations. Reweighting the training samples has aroused as an effective mitigation strategy to these problems. In this paper, we propose a novel and coherent scheme for kernel-reweighted regression by reparametrizing the sample weights using a doubly non-negative matrix. When the weighting matrix is confined in an uncertainty set using either the log-determinant divergence or the Bures-Wasserstein distance, we show that the adversarially reweighted estimate can be solved efficiently using first-order methods. Numerical experiments show that our reweighting strategy delivers promising results on numerous datasets.

[Learned Robust PCA: A Scalable Deep Unfolding Approach for High-Dimensional Outlier Detection](#)

- HanQin Cai, Jialin Liu, Wotao Yin
- abstract@[open-review](#): Robust principal component analysis (RPCA) is a critical tool in modern machine learning, which detects outliers in the task of low-rank matrix reconstruction. In this paper, we propose a scalable and learnable non-convex approach for high-dimensional RPCA problems, which we call Learned Robust PCA (LRPCA). LRPCA is highly efficient, and its free parameters can be effectively learned to optimize via deep unfolding. Moreover, we extend deep unfolding from finite iterations to infinite iterations via a novel feedforward-recurrent-mixed neural network model. We establish the recovery guarantee of LRPCA under mild assumptions for RPCA. Numerical experiments show that LRPCA outperforms the state-of-the-art RPCA algorithms, such as ScaledGD and AltProj, on both synthetic datasets and real-world applications.

[Proxy-Normalizing Activations to Match Batch Normalization while Removing Batch Dependence](#)

- Antoine Labatie, Dominic Masters, Zach Eaton-Rosen, Carlo Luschi
- abstract@[open-review](#): We investigate the reasons for the performance degradation incurred with batch-independent normalization. We find that the prototypical techniques of layer normalization and instance normalization both induce the appearance of failure modes in the neural network's pre-activations: (i) layer normalization induces a collapse towards channel-wise constant functions; (ii) instance normalization induces a lack of variability in instance statistics, symptomatic of an alteration of the expressivity. To alleviate failure mode (i) without aggravating failure mode (ii), we introduce the technique "Proxy Normalization" that normalizes post-activations using a proxy distribution. When combined with layer normalization or group normalization, this batch-independent normalization emulates batch normalization's behavior and consistently matches or exceeds its performance.

[Dynamic Bottleneck for Robust Self-Supervised Exploration](#)

- Chenjia Bai, Lingxiao Wang, Lei Han, Animesh Garg, Jianye Hao, Peng Liu, Zhaoran Wang
- abstract@[open-review](#): Exploration methods based on pseudo-count of transitions or curiosity of dynamics have achieved promising results in solving reinforcement learning with sparse rewards. However, such methods are usually sensitive to environmental dynamics-irrelevant information, e.g., white-noise. To handle such dynamics-irrelevant information, we propose a Dynamic Bottleneck (DB) model, which attains a dynamics-relevant representation based on the information-bottleneck principle. Based on the DB model, we further propose DB-bonus, which encourages the agent to explore state-action pairs with high information gain. We establish theoretical connections between the proposed DB-bonus, the upper confidence bound (UCB) for linear case, and the visiting count for tabular case. We evaluate the proposed method on Atari suits with dynamics-irrelevant noises. Our experiments show that exploration with DB bonus outperforms several state-of-the-art exploration methods in noisy environments.

[ProTo: Program-Guided Transformer for Program-Guided Tasks](#)

- Zelin Zhao, Karan Samel, Binghong Chen, lee song
- abstract@[open-review](#): Programs, consisting of semantic and structural information, play an important role in the communication between humans and agents. Towards learning general program executors to unify perception, reasoning, and decision making, we formulate program-guided tasks which require learning to execute a given program on the observed task specification. Furthermore, we propose Program-Guided Transformers (ProTo), which integrates both semantic and structural guidance of a program by leveraging cross-attention and masked self-attention to pass messages between the specification and routines in the program. ProTo executes a program in a learned latent space and enjoys stronger representation ability than previous neural-symbolic approaches. We demonstrate that ProTo significantly outperforms the previous state-of-the-art methods on GQA visual reasoning and 2D Minecraft policy learning datasets. Additionally, ProTo demonstrates better generalization to unseen, complex, and human-written programs.

[An Efficient Transfer Learning Framework for Multiagent Reinforcement Learning](#)

- Tianpei Yang, Weixun Wang, Hongyao Tang, Jianye Hao, Zhaopeng Meng, Hangyu Mao, Dong Li, Wulong Liu, Yingfeng Chen, Yujing Hu, Changjie Fan, Chengwei Zhang
- abstract@[open-review](#): Transfer Learning has shown great potential to enhance single-agent Reinforcement Learning (RL) efficiency. Similarly, Multiagent RL (MARL) can also be accelerated if agents can share knowledge with each other. However, it remains a problem of how an agent should learn from other agents. In this paper, we propose a novel Multiagent Policy Transfer Framework (MAPTF) to improve MARL efficiency. MAPTF learns which agent's policy is the best to reuse for each agent and when to terminate it by modeling multiagent policy transfer as the option learning problem. Furthermore, in practice, the option module can only collect all agent's local experiences for update due to the partial observability of the environment.

While in this setting, each agent's experience may be inconsistent with each other, which may cause the inaccuracy and oscillation of the option-value's estimation. Therefore, we propose a novel option learning algorithm, the successor representation option learning to solve it by decoupling the environment dynamics from rewards and learning the option-value under each agent's preference. MAPTF can be easily combined with existing deep RL and MARL approaches, and experimental results show it significantly boosts the performance of existing methods in both discrete and continuous state spaces.

[Learning to Time-Decode in Spiking Neural Networks Through the Information Bottleneck](#)

- Nicolas Skatchkovsky, Osvaldo Simeone, Hyeryung Jang
- abstract@[open-review](#): One of the key challenges in training Spiking Neural Networks (SNNs) is that target outputs typically come in the form of natural signals, such as labels for classification or images for generative models, and need to be encoded into spikes. This is done by handcrafting target spiking signals, which in turn implicitly fixes the mechanisms used to decode spikes into natural signals, e.g., rate decoding. The arbitrary choice of target signals and decoding rule generally impairs the capacity of the SNN to encode and process information in the timing of spikes. To address this problem, this work introduces a hybrid variational autoencoder architecture, consisting of an encoding SNN and a decoding Artificial Neural Network (ANN). The role of the decoding ANN is to learn how to best convert the spiking signals output by the SNN into the target natural signal. A novel end-to-end learning rule is introduced that optimizes a directed information bottleneck training criterion via surrogate gradients. We demonstrate the applicability of the technique in an experimental settings on various tasks, including real-life datasets.

[NEO: Non Equilibrium Sampling on the Orbits of a Deterministic Transform](#)

- Achille Thin, Yazid Janati El Idrissi, Sylvain Le Corff, Charles Ollion, Eric Moulines, Arnaud Doucet, Alain Durmus, Christian X Robert
- abstract@[open-review](#): Sampling from a complex distribution π and approximating its intractable normalizing constant Z are challenging problems. In this paper, a novel family of importance samplers (IS) and Markov chain Monte Carlo (MCMC) samplers is derived. Given an invertible map T , these schemes combine (with weights) elements from the forward and backward Orbits through points sampled from a proposal distribution ρ . The map T does not leave the target π invariant, hence the name NEO, standing for Non-Equilibrium Orbits. NEO-IS provides unbiased estimators of the normalizing constant and self-normalized IS estimators of expectations under π while NEO-MCMC combines multiple NEO-IS estimates of the normalizing constant and an iterated sampling-importance resampling mechanism to sample from π . For T chosen as a discrete-time integrator of a conformal Hamiltonian system, NEO-IS achieves state-of-the art performance on difficult benchmarks and NEO-MCMC is able to explore highly multimodal targets. Additionally, we provide detailed theoretical results for both methods. In particular, we show that NEO-MCMC is uniformly geometrically ergodic and establish explicit mixing time estimates under mild conditions.

[Relaxing Local Robustness](#)

- Klas Leino, Matt Fredrikson
- abstract@[open-review](#): Certifiable local robustness, which rigorously precludes small-norm adversarial examples, has received significant attention as a means of addressing security concerns in deep learning. However, for some classification problems, local robustness is not a natural objective, even in the presence of adversaries; for example, if an image contains two classes of subjects, the correct label for the image may be considered arbitrary between the two, and thus enforcing strict separation between them is unnecessary. In this work, we introduce two relaxed safety properties for classifiers that address this observation: (1) relaxed top-k robustness, which serves as the analogue of top-k accuracy; and (2) affinity robustness, which specifies which sets of labels must be separated by a robustness margin, and which can be ϵ -close in ℓ_p space. We show how to construct models that can be efficiently certified against each relaxed robustness property, and trained with very little overhead relative to standard gradient descent. Finally, we demonstrate experimentally that these relaxed variants of robustness are well-suited to several significant classification problems, leading to lower rejection rates and higher certified accuracies than can be obtained when certifying "standard" local robustness.

[Tuning Large Neural Networks via Zero-Shot Hyperparameter Transfer](#)

- Ge Yang, Edward Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, Jianfeng Gao
- abstract@[open-review](#): Hyperparameter (HP) tuning in deep learning is an expensive process, prohibitively so for neural networks (NNs) with billions of parameters. We show that, in the recently discovered Maximal Update Parametrization (μP), many optimal HPs remain stable even as model size changes. This leads to a new HP tuning paradigm we call μ Transfer: parametrize the target model in μP , tune the HP indirectly on a smaller model, and *zero-shot transfer* them to the full-sized model, i.e., without directly tuning the latter at all. We verify μ Transfer on Transformer and ResNet. For example, 1) by transferring pretraining HPs from a model of 13M parameters, we outperform published numbers of BERT-large (350M parameters), with a total tuning cost equivalent to pretraining BERT-large once; 2) by transferring from 40M parameters, we outperform published numbers of the 6.7B GPT-3 model, with tuning cost only 7% of total pretraining cost. A Pytorch implementation of our technique can be found at github.com/microsoft/mup. See arxiv.org for the full, up-to-date version of this work.

[Statistical Regeneration Guarantees of the Wasserstein Autoencoder with Latent Space Consistency](#)

- Anish Chakrabarty, Swagatam Das
- abstract@[open-review](#): The introduction of Variational Autoencoders (VAE) has been marked as a breakthrough in the history of representation learning models. Besides having several accolades of its own, VAE has successfully flagged off a series of inventions in the form of its immediate successors. Wasserstein Autoencoder (WAE), being an heir to that realm carries with it all of the goodness and heightened generative promises, matching even the generative adversarial networks (GANs). Needless to say, recent years have witnessed a remarkable resurgence in statistical analyses of the GANs. Similar examinations for Autoencoders however, despite their diverse applicability and notable empirical performance, remain largely absent. To close this gap, in this paper, we investigate the statistical properties of WAE. Firstly, we provide statistical guarantees that WAE achieves the target distribution in the latent space, utilizing the Vapnik–Chervonenkis (VC) theory. The main result, consequently ensures the regeneration of the input distribution, harnessing the potential offered by Optimal Transport of measures under the Wasserstein metric. This study, in turn, hints at the class of distributions WAE can reconstruct after suffering a compression in the form of a latent law.

[Leveraging the Inductive Bias of Large Language Models for Abstract Textual Reasoning](#)

- Christopher Rytting, David Wingate
- abstract@[open-review](#): Large natural language models (LMs) (such as GPT-3 or T5) demonstrate impressive abilities across a range of general NLP tasks. Here, we show that the knowledge embedded in such models provides a useful inductive bias, not just on traditional NLP tasks, but also in the nontraditional task of training a symbolic reasoning engine. We observe that these engines learn quickly and generalize in a natural way that reflects human intuition. For example, training such a system to model block-stacking might naturally generalize to stacking other types of objects because of structure in the real world that has been partially captured by the language describing it. We study several abstract textual reasoning tasks, such as object manipulation and navigation, and demonstrate multiple types of generalization to novel scenarios and the symbols that comprise them. We also demonstrate the surprising utility of $\text{compositional learning}$, where a learner dedicated to mastering a complicated task gains an advantage by training on relevant simpler tasks instead of jumping straight to the complicated task.

[Differentiable Simulation of Soft Multi-body Systems](#)

- Yiling Qiao, Junbang Liang, Vladlen Koltun, Ming Lin
- abstract@[open-review](#): We present a method for differentiable simulation of soft articulated bodies. Our work enables the integration of differentiable physical dynamics into gradient-based pipelines. We develop a top-down matrix assembly algorithm within Projective Dynamics and derive a generalized dry friction model for soft continuum using a new matrix splitting strategy. We derive a differentiable control framework for soft articulated bodies driven by muscles, joint torques, or pneumatic tubes. The experiments demonstrate that our designs make soft body simulation more stable and realistic compared to other frameworks. Our method accelerates the solution of system identification problems by more than an order of magnitude, and enables efficient gradient-based learning of motion control with soft robots.

[Good Classification Measures and How to Find Them](#)

- Martijn Gösgens, Anton Zhiyanov, Aleksey Tikhonov, Liudmila Prokhorenkova
- abstract@[open-review](#): Several performance measures can be used for evaluating classification results: accuracy, F-measure, and many others. Can we say that some of them are better than others, or, ideally, choose one measure that is best in all situations? To answer this question, we conduct a systematic analysis of classification performance measures: we formally define a list of desirable properties and theoretically analyze which measures satisfy which properties. We also prove an impossibility theorem: some desirable properties cannot be simultaneously satisfied. Finally, we propose a new family of measures satisfying all desirable properties except one. This family includes the Matthews Correlation Coefficient and a so-called Symmetric Balanced Accuracy that was not previously used in classification literature. We believe that our systematic approach gives an important tool to practitioners for adequately evaluating classification results.

[Distilling Robust and Non-Robust Features in Adversarial Examples by Information Bottleneck](#)

- Junho Kim, Byung-Kwan Lee, Yong Man Ro
- abstract@[open-review](#): Adversarial examples, generated by carefully crafted perturbation, have attracted considerable attention in research fields. Recent works have argued that the existence of the robust and non-robust features is a primary cause of the adversarial examples, and investigated their internal interactions in the feature space. In this paper, we propose a way of explicitly distilling feature representation into the robust and non-robust features, using Information Bottleneck. Specifically, we inject noise variation to each feature unit and evaluate the information flow in the feature representation to dichotomize feature units either robust or non-robust, based on the noise variation magnitude. Through comprehensive experiments, we demonstrate that the distilled features are highly correlated with adversarial prediction, and they have human-perceptible semantic information by themselves. Furthermore, we present an attack mechanism intensifying the gradient of non-robust features that is directly related to the model prediction, and validate its effectiveness of breaking model robustness.

[Vector-valued Gaussian Processes on Riemannian Manifolds via Gauge Independent Projected Kernels](#)

- Michael Hutchinson, Alexander Terenin, Viacheslav Borovitskiy, So Takao, Yee Teh, Marc Deisenroth
- abstract@[open-review](#): Gaussian processes are machine learning models capable of learning unknown functions in a way that represents uncertainty, thereby facilitating construction of optimal decision-making systems. Motivated by a desire to deploy Gaussian processes in novel areas of science, a rapidly-growing line of research has focused on constructively extending these models to handle non-Euclidean domains, including Riemannian manifolds, such as spheres and tori. We propose techniques that generalize this class to model vector fields on Riemannian manifolds, which are important in a number of application areas in the physical sciences. To do so, we present a general recipe for constructing gauge independent kernels, which induce Gaussian vector fields, i.e. vector-valued Gaussian processes coherent with geometry, from scalar-valued Riemannian kernels. We extend standard Gaussian process training methods, such as variational inference, to this setting. This enables vector-valued Gaussian processes on Riemannian manifolds to be trained using standard methods and makes them accessible to machine learning practitioners.

[On the Representation Power of Set Pooling Networks](#)

- Christian Bueno, Alan Hylton
- abstract@[open-review](#): Point clouds and sets are input data-types which pose unique problems to deep learning. Since sets can have variable cardinality and are unchanged by permutation, the input space for these problems naturally form infinite-dimensional non-Euclidean spaces. Despite these mathematical difficulties, PointNet (Qi et al. 2017) and Deep Sets (Zaheer et al. 2017) introduced foundational neural network architectures to address these problems. In this paper we present a unified framework to study the expressive power of such networks as well as their extensions beyond point clouds (partially addressing a conjecture on the extendibility of DeepSets along the way). To this end, we demonstrate the crucial role that the Hausdorff and Wasserstein metrics play and prove new cardinality-agnostic universality results to characterize exactly which functions can be approximated by these models. In particular, these results imply that PointNet generally cannot approximate averages of continuous functions over sets (e.g. center-of-mass or higher moments) implying that DeepSets is strictly more expressive than PointNet in the constant cardinality setting. Moreover, we obtain explicit lower-bounds on the approximation error and present a simple method to produce arbitrarily many examples of this failure-mode. Counterintuitively, we also prove that in the unbounded cardinality setting that any function which can be uniformly approximated by both PointNet and normalized-DeepSets must be constant. Finally, we also prove theorems on the Lipschitz properties of PointNet and normalized-DeepSets which shed insight into exploitable inductive bias in these networks.

[Learning Policies with Zero or Bounded Constraint Violation for Constrained MDPs](#)

- Tao Liu, Ruida Zhou, Dileep Kalathil, Panganamala Kumar, Chao Tian
- abstract@[open-review](#): We address the issue of safety in reinforcement learning. We pose the problem in an episodic framework of a constrained Markov decision process. Existing results have shown that it is possible to achieve a reward regret of $\tilde{\mathcal{O}}(\sqrt{K})$ while allowing an $\tilde{\mathcal{O}}(\sqrt{K})$ constraint violation in K episodes. A critical question that arises is whether it is possible to keep the constraint violation even smaller. We show that when a strictly safe policy is known, then one can confine the system to zero constraint violation with arbitrarily high probability while keeping the reward regret of order $\tilde{\mathcal{O}}(\sqrt{K})$. The algorithm which does so employs the principle of optimistic pessimism in the face of uncertainty to achieve safe exploration. When no strictly safe policy is known, though one is known to exist, then it is possible to restrict the system to bounded constraint violation with arbitrarily high probability. This is shown to be realized by a primal-dual algorithm with an optimistic primal estimate and a pessimistic dual update.

[A Prototype-Oriented Framework for Unsupervised Domain Adaptation](#)

- Korawat Tanwisuth, Xinjie Fan, Huangjie Zheng, Shujian Zhang, Hao Zhang, Bo Chen, Mingyuan Zhou
- abstract@[open-review](#): Existing methods for unsupervised domain adaptation often rely on minimizing some statistical distance between the source and target samples in the latent space. To avoid the sampling variability, class imbalance, and data-privacy concerns that often plague these methods, we instead provide a memory and computation-efficient probabilistic framework to extract class prototypes and align the target features with them. We demonstrate the general applicability of our method on a wide range of scenarios, including single-source, multi-source, class-imbalance, and source-

private domain adaptation. Requiring no additional model parameters and having a moderate increase in computation over the source model alone, the proposed method achieves competitive performance with state-of-the-art methods.

[Mining the Benefits of Two-stage and One-stage HOI Detection](#)

- Aixi Zhang, Yue Liao, Si Liu, Miao Lu, Yongliang Wang, Chen Gao, XIAOBO LI
- abstract@[open-review](#): Two-stage methods have dominated Human-Object Interaction~(HOI) detection for several years. Recently, one-stage HOI detection methods have become popular. In this paper, we aim to explore the essential pros and cons of two-stage and one-stage methods. With this as the goal, we find that conventional two-stage methods mainly suffer from positioning positive interactive human-object pairs, while one-stage methods are challenging to make an appropriate trade-off on multi-task learning, \emph{i.e.}, object detection, and interaction classification. Therefore, a core problem is how to take the essence and discard the dregs from the conventional two types of methods. To this end, we propose a novel one-stage framework with disentangling human-object detection and interaction classification in a cascade manner. In detail, we first design a human-object pair generator based on a state-of-the-art one-stage HOI detector by removing the interaction classification module or head and then design a relatively isolated interaction classifier to classify each human-object pair. Two cascade decoders in our proposed framework can focus on one specific task, detection or interaction classification. In terms of the specific implementation, we adopt a transformer-based HOI detector as our base model. The newly introduced disentangling paradigm outperforms existing methods by a large margin, with a significant relative mAP gain of 9.32% on HICO-Det. The source codes are available at <https://github.com/YueLiao/CDN>.

[Discerning Decision-Making Process of Deep Neural Networks with Hierarchical Voting Transformation](#)

- Ying Sun, Hengshu Zhu, Chuan Qin, Fuzhen Zhuang, Qing He, Hui Xiong
- abstract@[open-review](#): Neural network based deep learning techniques have shown great success for numerous applications. While it is expected to understand their intrinsic decision-making processes, these deep neural networks often work in a black-box way. To this end, in this paper, we aim to discern the decision-making processes of neural networks through a hierarchical voting strategy by developing an explainable deep learning model, namely Voting Transformation-based Explainable Neural Network (VOTEN). Specifically, instead of relying on massive feature combinations, VOTEN creatively models expressive single-valued voting functions between explicitly modeled latent concepts to achieve high fitting ability. Along this line, we first theoretically analyze the major components of VOTEN and prove the relationship and advantages of VOTEN compared with Multi-Layer Perceptron (MLP), the basic structure of deep neural networks. Moreover, we design efficient algorithms to improve the model usability by explicitly showing the decision processes of VOTEN. Finally, extensive experiments on multiple real-world datasets clearly validate the performances and explainability of VOTEN.

[Risk-averse Heteroscedastic Bayesian Optimization](#)

- Anastasia Makarova, Ilnura Usmanova, Ilija Bogunovic, Andreas Krause
- abstract@[open-review](#): Many black-box optimization tasks arising in high-stakes applications require risk-averse decisions. The standard Bayesian optimization (BO) paradigm, however, optimizes the expected value only. We generalize BO to trade mean and input-dependent variance of the objective, both of which we assume to be unknown a priori. In particular, we propose a novel risk-averse heteroscedastic Bayesian optimization algorithm (RAHBO) that aims to identify a solution with high return and low noise variance, while learning the noise distribution on the fly. To this end, we model both expectation and variance as (unknown) RKHS functions, and propose a novel risk-aware acquisition function. We bound the regret for our approach and provide a robust rule to report the final decision point for applications where only a single solution must be identified. We demonstrate the effectiveness of RAHBO on synthetic benchmark functions and hyperparameter tuning tasks.

[Invertible DenseNets with Concatenated LipSwish](#)

- Yura Perugachi-Diaz, Jakub Tomczak, Sandjai Bhulai
- abstract@[open-review](#): We introduce Invertible Dense Networks (i-DenseNets), a more parameter efficient extension of Residual Flows. The method relies on an analysis of the Lipschitz continuity of the concatenation in DenseNets, where we enforce invertibility of the network by satisfying the Lipschitz constant. Furthermore, we propose a learnable weighted concatenation, which not only improves the model performance but also indicates the importance of the concatenated weighted representation. Additionally, we introduce the Concatenated LipSwish as activation function, for which we show how to enforce the Lipschitz condition and which boosts performance. The new architecture, i-DenseNet, out-performs Residual Flow and other flow-based models on density estimation evaluated in bits per dimension, where we utilize an equal parameter budget. Moreover, we show that the proposed model out-performs Residual Flows when trained as a hybrid model where the model is both a generative and a discriminative model.

[Topological Detection of Trojaned Neural Networks](#)

- Songzhu Zheng, Yikai Zhang, Hubert Wagner, Mayank Goswami, Chao Chen
- abstract@[open-review](#): Deep neural networks are known to have security issues. One particular threat is the Trojan attack. It occurs when the attackers stealthily manipulate the model's behavior through Trojaned training samples, which can later be exploited. Guided by basic neuroscientific principles, we discover subtle -- yet critical -- structural deviation characterizing Trojaned models. In our analysis we use topological tools. They allow us to model high-order dependencies in the networks, robustly compare different networks, and localize structural abnormalities. One interesting observation is that Trojaned models develop short-cuts from shallow to deep layers. Inspired by these observations, we devise a strategy for robust detection of Trojaned models. Compared to standard baselines it displays better performance on multiple benchmarks.

[Provably Strict Generalisation Benefit for Invariance in Kernel Methods](#)

- Bryn Elesedy
- abstract@[open-review](#): It is a commonly held belief that enforcing invariance improves generalisation. Although this approach enjoys widespread popularity, it is only very recently that a rigorous theoretical demonstration of this benefit has been established. In this work we build on the function space perspective of Elesedy and Zaidi [8] to derive a strictly non-zero generalisation benefit of incorporating invariance in kernel ridge regression when the target is invariant to the action of a compact group. We study invariance enforced by feature averaging and find that generalisation is governed by a notion of effective dimension that arises from the interplay between the kernel and the group. In building towards this result, we find that the action of the group induces an orthogonal decomposition of both the reproducing kernel Hilbert space and its kernel, which may be of interest in its own right.

[Formalizing the Generalization-Forgetting Trade-off in Continual Learning](#)

- Krishnan Raghavan, Prasanna Balaprakash
- abstract@[open-review](#): We formulate the continual learning (CL) problem via dynamic programming and model the trade-off between catastrophic forgetting and generalization as a two-player sequential game. In this approach, player 1 maximizes the cost due to lack of generalization whereas player 2 minimizes the cost due to catastrophic forgetting. We show theoretically that a balance point between the two players exists for each task and that this point is stable (once the balance is achieved, the two players stay at the balance point). Next, we introduce balanced continual learning (BCL), which is designed to attain balance between generalization and forgetting and empirically demonstrate that BCL is comparable to or better than the state of the art.

Risk-Aware Transfer in Reinforcement Learning using Successor Features

- Michael Gimelfarb, Andre Barreto, Scott Sanner, Chi-Guhn Lee
- abstract@[open-review](#): Sample efficiency and risk-awareness are central to the development of practical reinforcement learning (RL) for complex decision-making. The former can be addressed by transfer learning, while the latter by optimizing some utility function of the return. However, the problem of transferring skills in a risk-aware manner is not well-understood. In this paper, we address the problem of transferring policies between tasks in a common domain that differ only in their reward functions, in which risk is measured by the variance of reward streams. Our approach begins by extending the idea of generalized policy improvement to maximize entropic utilities, thus extending the dynamic programming's policy improvement operation to sets of policies \emph{and} levels of risk-aversion. Next, we extend the idea of successor features (SF), a value function representation that decouples the environment dynamics from the rewards, to capture the variance of returns. Our resulting risk-aware successor features (RaSF) integrate seamlessly within the RL framework, inherit the superior task generalization ability of SFs, while incorporating risk into the decision-making. Experiments on a discrete navigation domain and control of a simulated robotic arm demonstrate the ability of RaSFs to outperform alternative methods including SFs, when taking the risk of the learned policies into account.

Causal Inference for Event Pairs in Multivariate Point Processes

- Tian Gao, Dharmashankar Subramanian, Debarun Bhattacharjya, Xiao Shou, Nicholas Mattei, Kristin P Bennett
- abstract@[open-review](#): Causal inference and discovery from observational data has been extensively studied across multiple fields. However, most prior work has focused on independent and identically distributed (i.i.d.) data. In this paper, we propose a formalization for causal inference between pairs of event variables in multivariate recurrent event streams by extending Rubin's framework for the average treatment effect (ATE) and propensity scores to multivariate point processes. Analogous to a joint probability distribution representing i.i.d. data, a multivariate point process represents data involving asynchronous and irregularly spaced occurrences of various types of events over a common timeline. We theoretically justify our point process causal framework and show how to obtain unbiased estimates of the proposed measure. We conduct an experimental investigation using synthetic and real-world event datasets, where our proposed causal inference framework is shown to exhibit superior performance against a set of baseline pairwise causal association scores.

Evaluating model performance under worst-case subpopulations

- Mike Li, Hongseok Namkoong, Shangzhou Xia
- abstract@[open-review](#): The performance of ML models degrades when the training population is different from that seen under operation. Towards assessing distributional robustness, we study the worst-case performance of a model over all subpopulations of a given size, defined with respect to core attributes Z . This notion of robustness can consider arbitrary (continuous) attributes Z , and automatically accounts for complex intersectionality in disadvantaged groups. We develop a scalable yet principled two-stage estimation procedure that can evaluate the robustness of state-of-the-art models. We prove that our procedure enjoys several finite-sample convergence guarantees, including dimension-free convergence. Instead of overly conservative notions based on Rademacher complexities, our evaluation error depends on the dimension of Z only through the out-of-sample error in estimating the performance conditional on Z . On real datasets, we demonstrate that our method certifies the robustness of a model and prevents deployment of unreliable models.

Privately Publishable Per-instance Privacy

- Rachel Redberg, Yu-Xiang Wang
- abstract@[open-review](#): We consider how to privately share the personalized privacy losses incurred by objective perturbation, using per-instance differential privacy (pDP). Standard differential privacy (DP) gives us a worst-case bound that might be orders of magnitude larger than the privacy loss to a particular individual relative to a fixed dataset. The pDP framework provides a more fine-grained analysis of the privacy guarantee to a target individual, but the per-instance privacy loss itself might be a function of sensitive data. In this paper, we analyze the per-instance privacy loss of releasing a private empirical risk minimizer learned via objective perturbation, and propose a group of methods to privately and accurately publish the pDP losses at little to no additional privacy cost.

Understanding the Limits of Unsupervised Domain Adaptation via Data Poisoning

- Akshay Mehra, Bhavya Kailkhura, Pin-Yu Chen, Jihun Hamm
- abstract@[open-review](#): Unsupervised domain adaptation (UDA) enables cross-domain learning without target domain labels by transferring knowledge from a labeled source domain whose distribution differs from that of the target. However, UDA is not always successful and several accounts of 'negative transfer' have been reported in the literature. In this work, we prove a simple lower bound on the target domain error that complements the existing upper bound. Our bound shows the insufficiency of minimizing source domain error and marginal distribution mismatch for a guaranteed reduction in the target domain error, due to the possible increase of induced labeling function mismatch. This insufficiency is further illustrated through simple distributions for which the same UDA approach succeeds, fails, and may succeed or fail with an equal chance. Motivated from this, we propose novel data poisoning attacks to fool UDA methods into learning representations that produce large target domain errors. We evaluate the effect of these attacks on popular UDA methods using benchmark datasets where they have been previously shown to be successful. Our results show that poisoning can significantly decrease the target domain accuracy, dropping it to almost 0% in some cases, with the addition of only 10% poisoned data in the source domain. The failure of these UDA methods demonstrates their limitations at guaranteeing cross-domain generalization consistent with our lower bound. Thus, evaluating UDA methods in adversarial settings such as data poisoning provides a better sense of their robustness to data distributions unfavorable for UDA.

Coresets for Clustering with Missing Values

- Vladimir Braverman, Shaofeng Jiang, Robert Krauthgamer, Xuan Wu
- abstract@[open-review](#): We provide the first coresnet for clustering points in \mathbb{R}^d that have multiple missing values (coordinates). Previous coresnet constructions only allow one missing coordinate. The challenge in this setting is that objective functions, like kMeans, are evaluated only on the set of available (non-missing) coordinates, which varies across points. Recall that an ϵ -coresnet of a large dataset is a small proxy, usually a reweighted subset of points, that $(1+\epsilon)$ -approximates the clustering objective for every possible center set. Our coresnets for k -Means and k -Median clustering have size $O(\min(j,k)) \cdot (\epsilon^{-1} d \log n)^2$, where n is the number of data points, d is the dimension and j is the maximum number of missing coordinates for each data point. We further design an algorithm to construct these coresnets in near-linear time, and consequently improve a recent quadratic-time PTAS for k -Means with missing values [Eiben et al., SODA 2021] to near-linear time. We validate our coresnet construction, which is based on importance sampling and is easy to implement, on various real data sets. Our coresnet exhibits a flexible tradeoff between coresnet size and accuracy, and generally outperforms the uniform-sampling baseline. Furthermore, it significantly speeds up a Lloyd's-style heuristic for k -Means with missing values.

Boosting with Multiple Sources

- Corinna Cortes, Mehryar Mohri, Dmitry Storcheus, Ananda Theertha Suresh

- abstract@[open-review](#): We study the problem of learning accurate ensemble predictors, in particular boosting, in the presence of multiple source domains. We show that the standard convex combination ensembles in general cannot succeed in this scenario and adopt instead a domain-weighted combination. We introduce and analyze a new boosting algorithm, MULTIBOOST, for this scenario and show that it benefits from favorable theoretical guarantees. We also report the results of several experiments with our algorithm demonstrating that it outperforms natural baselines on multi-source text-based, image-based and tabular data. We further present an extension of our algorithm to the federated learning scenario and report favorable experimental results for that setting as well. Additionally, we describe in detail an extension of our algorithm to the multi-class setting, MCMULTIBOOST, for which we also report experimental results.

[Dynamic Neural Representational Decoders for High-Resolution Semantic Segmentation](#)

- Bowen Zhang, Yifan Liu, Zhi Tian, Chunhua Shen
- abstract@[open-review](#): Semantic segmentation requires per-pixel prediction for a given image. Typically, the output resolution of a segmentation network is severely reduced due to the downsampling operations in the CNN backbone. Most previous methods employ upsampling decoders to recover the spatial resolution. Various decoders were designed in the literature. Here, we propose a novel decoder, termed dynamic neural representational decoder (NRD), which is simple yet significantly more efficient. As each location on the encoder's output corresponds to a local patch of the semantic labels, in this work, we represent these local patches of labels with compact neural networks. This neural representation enables our decoder to leverage the smoothness prior in the semantic label space, and thus makes our decoder more efficient. Furthermore, these neural representations are dynamically generated and conditioned on the outputs of the encoder networks. The desired semantic labels can be efficiently decoded from the neural representations, resulting in high-resolution semantic segmentation predictions. We empirically show that our proposed decoder can outperform the decoder in DeeplabV3+ with only $\sim 30\%$ computational complexity, and achieve competitive performance with the methods using dilated encoders with only $\sim 15\%$ computation. Experiments on Cityscapes, ADE20K, and Pascal Context demonstrate the effectiveness and efficiency of our proposed method.

[Dense Keypoints via Multiview Supervision](#)

- Zhixuan Yu, Haozheng Yu, Long Sha, Sujoy Ganguly, Hyun Soo Park
- abstract@[open-review](#): This paper presents a new end-to-end semi-supervised framework to learn a dense keypoint detector using unlabeled multiview images. A key challenge lies in finding the exact correspondences between the dense keypoints in multiple views since the inverse of the keypoint mapping can be neither analytically derived nor differentiated. This limits applying existing multiview supervision approaches used to learn sparse keypoints that rely on the exact correspondences. To address this challenge, we derive a new probabilistic epipolar constraint that encodes the two desired properties. (1) Soft correspondence: we define a matchability, which measures a likelihood of a point matching to the other image's corresponding point, thus relaxing the requirement of the exact correspondences. (2) Geometric consistency: every point in the continuous correspondence fields must satisfy the multiview consistency collectively. We formulate a probabilistic epipolar constraint using a weighted average of epipolar errors through the matchability thereby generalizing the point-to-point geometric error to the field-to-field geometric error. This generalization facilitates learning a geometrically coherent dense keypoint detection model by utilizing a large number of unlabeled multiview images. Additionally, to prevent degenerative cases, we employ a distillation-based regularization by using a pretrained model. Finally, we design a new neural network architecture, made of twin networks, that effectively minimizes the probabilistic epipolar errors of all possible correspondences between two view images by building affinity matrices. Our method shows superior performance compared to existing methods, including non-differentiable bootstrapping in terms of keypoint accuracy, multiview consistency, and 3D reconstruction accuracy.

[Scatterbrain: Unifying Sparse and Low-rank Attention](#)

- Beidi Chen, Tri Dao, Eric Winsor, Zhao Song, Atri Rudra, Christopher Ré
- abstract@[open-review](#): Recent advances in efficient Transformers have exploited either the sparsity or low-rank properties of attention matrices to reduce the computational and memory bottlenecks of modeling long sequences. However, it is still challenging to balance the trade-off between model quality and efficiency to perform a one-size-fits-all approximation for different tasks. To better understand this trade-off, we observe that sparse and low-rank approximations excel in different regimes, determined by the softmax temperature in attention, and sparse + low-rank can outperform each individually. Inspired by the classical robust-PCA algorithm for sparse and low-rank decomposition, we propose Scatterbrain, a novel way to unify sparse (via locality sensitive hashing) and low-rank (via kernel feature map) attention for accurate and efficient approximation. The estimation is unbiased with provably low error. We empirically show that Scatterbrain can achieve $\sim 2.1 \times$ lower error than baselines when serving as a drop-in replacement in BigGAN image generation and pre-trained T2T-ViT. On a pre-trained T2T Vision transformer, even without fine-tuning, Scatterbrain can reduce $\sim 98\%$ of attention memory at the cost of only $\sim 1\%$ drop in accuracy. We demonstrate Scatterbrain for end-to-end training with up to $\sim 4\%$ points better perplexity and 5 points better average accuracy than sparse or low-rank efficient transformers on language modeling and long-range-area tasks.

[PTR: A Benchmark for Part-based Conceptual, Relational, and Physical Reasoning](#)

- Yining Hong, Li Yi, Josh Tenenbaum, Antonio Torralba, Chuang Gan
- abstract@[open-review](#): A critical aspect of human visual perception is the ability to parse visual scenes into individual objects and further into object parts, forming part-whole hierarchies. Such composite structures could induce a rich set of semantic concepts and relations, thus playing an important role in the interpretation and organization of visual signals as well as for the generalization of visual perception and reasoning. However, existing visual reasoning benchmarks mostly focus on objects rather than parts. Visual reasoning based on the full part-whole hierarchy is much more challenging than object-centric reasoning due to finer-grained concepts, richer geometry relations, and more complex physics. Therefore, to better serve for part-based conceptual, relational and physical reasoning, we introduce a new large-scale diagnostic visual reasoning dataset named PTR. PTR contains around 80k RGBD synthetic images with ground truth object and part level annotations regarding semantic instance segmentation, color attributes, spatial and geometric relationships, and certain physical properties such as stability. These images are paired with 800k machine-generated questions covering various types of reasoning types, making them a good testbed for visual reasoning models. We examine several state-of-the-art visual reasoning models on this dataset and observe that they still make many surprising mistakes in situations where humans can easily infer the correct answer. We believe this dataset will open up new opportunities for part-based reasoning. PTR dataset and baseline models are publicly available.

[Property-Aware Relation Networks for Few-Shot Molecular Property Prediction](#)

- Yaqing Wang, Abulikemu Abuduweili, Quanming Yao, Dejing Dou
- abstract@[open-review](#): Molecular property prediction plays a fundamental role in drug discovery to identify candidate molecules with target properties. However, molecular property prediction is essentially a few-shot problem, which makes it hard to use regular machine learning models. In this paper, we propose Property-Aware Relation networks (PAR) to handle this problem. In comparison to existing works, we leverage the fact that both relevant substructures and relationships among molecules change across different molecular properties. We first introduce a property-aware embedding function to transform the generic molecular embeddings to substructure-aware space relevant to the target property. Further, we design an adaptive relation graph learning module to jointly estimate molecular relation graph and refine molecular embeddings w.r.t. the target property, such that the limited labels can be effectively propagated among similar molecules. We adopt a meta-learning strategy where the parameters are selectively updated within tasks in order to model generic and property-aware knowledge separately. Extensive experiments on benchmark molecular property prediction datasets show that PAR consistently outperforms existing methods and can obtain property-aware molecular embeddings and model molecular relation graph properly.

[Differentially Private Learning with Adaptive Clipping](#)

- Galen Andrew, Om Thakkar, Brendan McMahan, Swaroop Ramaswamy
- abstract@[open-review](#): Existing approaches for training neural networks with user-level differential privacy (e.g., DP Federated Averaging) in federated learning (FL) settings involve bounding the contribution of each user's model update by $\{\backslash em\}$ clipping it to some constant value. However there is no good $\{\backslash em\}$ a priori setting of the clipping norm across tasks and learning settings: the update norm distribution depends on the model architecture and loss, the amount of data on each device, the client learning rate, and possibly various other parameters. We propose a method wherein instead of a fixed clipping norm, one clips to a value at a specified quantile of the update norm distribution, where the value at the quantile is itself estimated online, with differential privacy. The method tracks the quantile closely, uses a negligible amount of privacy budget, is compatible with other federated learning technologies such as compression and secure aggregation, and has a straightforward joint DP analysis with DP-FedAvg. Experiments demonstrate that adaptive clipping to the median update norm works well across a range of federated learning tasks, eliminating the need to tune any clipping hyperparameter.

[Can Less be More? When Increasing-to-Balancing Label Noise Rates Considered Beneficial](#)

- Yang Liu, Jialu Wang
- abstract@[open-review](#): In this paper, we answer the question of when inserting label noise (less informative labels) can instead return us more accurate and fair models. We are primarily inspired by three observations: 1) In contrast to reducing label noise rates, increasing the noise rates is easy to implement; 2) Increasing a certain class of instances' label noise to balance the noise rates (increasing-to-balancing) results in an easier learning problem; 3) Increasing-to-balancing improves fairness guarantees against label bias. In this paper, we first quantify the trade-offs introduced by increasing a certain group of instances' label noise rate w.r.t. the loss of label informativeness and the lowered learning difficulties. We analytically demonstrate when such an increase is beneficial, in terms of either improved generalization power or the fairness guarantees. Then we present a method to insert label noise properly for the task of learning with noisy labels, either without or with a fairness constraint. The primary technical challenge we face is due to the fact that we would not know which data instances are suffering from higher noise, and we would not have the ground truth labels to verify any possible hypothesis. We propose a detection method that informs us which group of labels might suffer from higher noise without using ground truth labels. We formally establish the effectiveness of the proposed solution and demonstrate it with extensive experiments.

[Projected GANs Converge Faster](#)

- Axel Sauer, Kashyap Chitta, Jens Müller, Andreas Geiger
- abstract@[open-review](#): Generative Adversarial Networks (GANs) produce high-quality images but are challenging to train. They need careful regularization, vast amounts of compute, and expensive hyper-parameter sweeps. We make significant headway on these issues by projecting generated and real samples into a fixed, pretrained feature space. Motivated by the finding that the discriminator cannot fully exploit features from deeper layers of the pretrained model, we propose a more effective strategy that mixes features across channels and resolutions. Our Projected GAN improves image quality, sample efficiency, and convergence speed. It is further compatible with resolutions of up to one Megapixel and advances the state-of-the-art Fréchet Inception Distance (FID) on twenty-two benchmark datasets. Importantly, Projected GANs match the previously lowest FIDs up to 40 times faster, cutting the wall-clock time from 5 days to less than 3 hours given the same computational resources.

[Generating High-Quality Explanations for Navigation in Partially-Revealed Environments](#)

- Gregory Stein
- abstract@[open-review](#): We present an approach for generating natural language explanations of high-level behavior of autonomous agents navigating in partially-revealed environments. Our counterfactual explanations communicate changes to interpretable statistics of the belief (e.g., the likelihood an exploratory action will reach the unseen goal) that are estimated from visual input via a deep neural network and used (via a Bellman equation variant) to inform planning far into the future. Additionally, our novel training procedure mimics explanation generation, allowing us to use planning performance as an objective measure of explanation quality. Simulated experiments validate that our explanations are both high quality and can be used in interventions to directly correct bad behavior; agents trained via our training-by-explaining procedure achieve 9.1% lower average cost than a non-learned baseline (12.7% after interventions) in environments derived from real-world floor plans.

[De-randomizing MCMC dynamics with the diffusion Stein operator](#)

- Zheyang Shen, Markus Heinonen, Samuel Kaski
- abstract@[open-review](#): Approximate Bayesian inference estimates descriptors of an intractable target distribution - in essence, an optimization problem within a family of distributions. For example, Langevin dynamics (LD) extracts asymptotically exact samples from a diffusion process because the time evolution of its marginal distributions constitutes a curve that minimizes the KL-divergence via steepest descent in the Wasserstein space. Parallel to LD, Stein variational gradient descent (SVGD) similarly minimizes the KL, albeit endowed with a novel Stein-Wasserstein distance, by deterministically transporting a set of particle samples, thus de-randomizes the stochastic diffusion process. We propose de-randomized kernel-based particle samplers to all diffusion-based samplers known as MCMC dynamics. Following previous work in interpreting MCMC dynamics, we equip the Stein-Wasserstein space with a fiber-Riemannian Poisson structure, with the capacity of characterizing a fiber-gradient Hamiltonian flow that simulates MCMC dynamics. Such dynamics discretizes into generalized SVGD (GSVGD), a Stein-type deterministic particle sampler, with particle updates coinciding with applying the diffusion Stein operator to a kernel function. We demonstrate empirically that GSVGD can de-randomize complex MCMC dynamics, which combine the advantages of auxiliary momentum variables and Riemannian structure, while maintaining the high sample quality from an interacting particle system.

[Sparsely Changing Latent States for Prediction and Planning in Partially Observable Domains](#)

- Christian Gumbtsch, Martin V. Butz, Georg Martius
- abstract@[open-review](#): A common approach to prediction and planning in partially observable domains is to use recurrent neural networks (RNNs), which ideally develop and maintain a latent memory about hidden, task-relevant factors. We hypothesize that many of these hidden factors in the physical world are constant over time, changing only sparsely. To study this hypothesis, we propose Gated \$L_0\$ Regularized Dynamics (GateL0RD), a novel recurrent architecture that incorporates the inductive bias to maintain stable, sparsely changing latent states. The bias is implemented by means of a novel internal gating function and a penalty on the L_0 norm of latent state changes. We demonstrate that GateL0RD can compete with or outperform state-of-the-art RNNs in a variety of partially observable prediction and control tasks. GateL0RD tends to encode the underlying generative factors of the environment, ignores spurious temporal dependencies, and generalizes better, improving sampling efficiency and overall performance in model-based planning and reinforcement learning tasks. Moreover, we show that the developing latent states can be easily interpreted, which is a step towards better explainability in RNNs.

[PreferenceNet: Encoding Human Preferences in Auction Design with Deep Learning](#)

- Neehar Peri, Michael Curry, Samuel Dooley, John Dickerson
- abstract@[open-review](#): The design of optimal auctions is a problem of interest in economics, game theory and computer science. Despite decades of effort, strategyproof, revenue-maximizing auction designs are still not known outside of restricted settings. However, recent methods using deep learning have shown some success in approximating optimal auctions, recovering several known solutions and outperforming strong baselines when optimal auctions are not known. In addition to maximizing revenue, auction mechanisms may also seek to encourage socially desirable constraints such as allocation fairness or diversity. However, these philosophical notions neither have standardization nor do they have widely accepted formal definitions. In

this paper, we propose PreferenceNet, an extension of existing neural-network-based auction mechanisms to encode constraints using (potentially human-provided) exemplars of desirable allocations. In addition, we introduce a new metric to evaluate an auction allocations' adherence to such socially desirable constraints and demonstrate that our proposed method is competitive with current state-of-the-art neural-network based auction designs. We validate our approach through human subject research and show that we are able to effectively capture real human preferences.

[Large-Scale Learning with Fourier Features and Tensor Decompositions](#)

- Frederiek Wesel, Kim Batselier
- abstract@[open-review](#): Random Fourier features provide a way to tackle large-scale machine learning problems with kernel methods. Their slow Monte Carlo convergence rate has motivated the research of deterministic Fourier features whose approximation error can decrease exponentially in the number of basis functions. However, due to their tensor product extension to multiple dimensions, these methods suffer heavily from the curse of dimensionality, limiting their applicability to one, two or three-dimensional scenarios. In our approach we overcome said curse of dimensionality by exploiting the tensor product structure of deterministic Fourier features, which enables us to represent the model parameters as a low-rank tensor decomposition. We derive a monotonically converging block coordinate descent algorithm with linear complexity in both the sample size and the dimensionality of the inputs for a regularized squared loss function, allowing to learn a parsimonious model in decomposed form using deterministic Fourier features. We demonstrate by means of numerical experiments how our low-rank tensor approach obtains the same performance of the corresponding nonparametric model, consistently outperforming random Fourier features.

[Hash Layers For Large Sparse Models](#)

- Stephen Roller, Sainbayar Sukhbaatar, arthur szlam, Jason Weston
- abstract@[open-review](#): We investigate the training of sparse layers that use different parameters for different inputs based on hashing in large Transformer models. Specifically, we modify the feedforward layer to hash to different sets of weights depending on the current token, over all tokens in the sequence. We show that this procedure either outperforms or is competitive with learning-to-route mixture-of-expert methods such as Switch Transformers and BASE Layers, while requiring no routing parameters or extra terms in the objective function such as a load balancing loss, and no sophisticated assignment algorithm. We study the performance of different hashing techniques, hash sizes and input features, and show that balanced and random hashes focused on the most local features work best, compared to either learning clusters or using longer-range context. We show our approach works well both on large language modeling and dialogue tasks, and on downstream fine-tuning tasks.

[Sliced Mutual Information: A Scalable Measure of Statistical Dependence](#)

- Ziv Goldfeld, Kristjan Greenewald
- abstract@[open-review](#): Mutual information (MI) is a fundamental measure of statistical dependence, with a myriad of applications to information theory, statistics, and machine learning. While it possesses many desirable structural properties, the estimation of high-dimensional MI from samples suffers from the curse of dimensionality. Motivated by statistical scalability to high dimensions, this paper proposes sliced MI (SMI) as a surrogate measure of dependence. SMI is defined as an average of MI terms between one-dimensional random projections. We show that it preserves many of the structural properties of classic MI, while gaining scalable computation and efficient estimation from samples. Furthermore, and in contrast to classic MI, SMI can grow as a result of deterministic transformations. This enables leveraging SMI for feature extraction by optimizing it over processing functions of raw data to identify useful representations thereof. Our theory is supported by numerical studies of independence testing and feature extraction, which demonstrate the potential gains SMI offers over classic MI for high-dimensional inference.

[Emergent Communication under Varying Sizes and Connectivities](#)

- Jooyeon Kim, Alice Oh
- abstract@[open-review](#): Recent advances in deep neural networks allowed artificial agents to derive their own emergent languages that promote interaction, coordination, and collaboration within a group. Just as we humans have succeeded in creating a shared language that allows us to interact within a large group, can the emergent communication within an artificial group converge to a shared, agreed language? This research provides an analytical study of the shared emergent language within the group communication settings of different sizes and connectivities. As the group size increases up to hundreds, agents start to speak dissimilar languages, but the rate at which they successfully communicate is maintained. We observe the emergence of different dialects when we restrict the group communication to have local connectivities only. Finally, we provide optimization results of group communication graphs when the number of agents one can communicate with is restricted or when we penalize communication between distant agent pairs. The optimized communication graphs show superior communication success rates compared to graphs with same number of links as well as the emergence of hub nodes and scale-free networks.

[Deep Bandits Show-Off: Simple and Efficient Exploration with Deep Networks](#)

- Rong Zhu, Mattia Rigotti
- abstract@[open-review](#): Designing efficient exploration is central to Reinforcement Learning due to the fundamental problem posed by the exploration-exploitation dilemma. Bayesian exploration strategies like Thompson Sampling resolve this trade-off in a principled way by modeling and updating the distribution of the parameters of the action-value function, the outcome model of the environment. However, this technique becomes infeasible for complex environments due to the computational intractability of maintaining probability distributions over parameters of outcome models of corresponding complexity. Moreover, the approximation techniques introduced to mitigate this issue typically result in poor exploration-exploitation trade-offs, as observed in the case of deep neural network models with approximate posterior methods that have been shown to underperform in the deep bandit scenario. In this paper we introduce Sample Average Uncertainty (SAU), a simple and efficient uncertainty measure for contextual bandits. While Bayesian approaches like Thompson Sampling estimate outcomes uncertainty indirectly by first quantifying the variability over the parameters of the outcome model, SAU is a frequentist approach that directly estimates the uncertainty of the outcomes based on the value predictions. Importantly, we show theoretically that the uncertainty measure estimated by SAU asymptotically matches the uncertainty provided by Thompson Sampling, as well as its regret bounds. Because of its simplicity SAU can be seamlessly applied to deep contextual bandits as a very scalable drop-in replacement for epsilon-greedy exploration. We confirm empirically our theory by showing that SAU-based exploration outperforms current state-of-the-art deep Bayesian bandit methods on several real-world datasets at modest computation cost, and make the code to reproduce our results available at \url{https://github.com/ibm/sau-explore}.

[Regret Minimization Experience Replay in Off-Policy Reinforcement Learning](#)

- Xu-Hui Liu, Zhenghai Xue, Jingcheng Pang, Shengyi Jiang, Feng Xu, Yang Yu
- abstract@[open-review](#): In reinforcement learning, experience replay stores past samples for further reuse. Prioritized sampling is a promising technique to better utilize these samples. Previous criteria of prioritization include TD error, recency and corrective feedback, which are mostly heuristically designed. In this work, we start from the regret minimization objective, and obtain an optimal prioritization strategy for Bellman update that can directly maximize the return of the policy. The theory suggests that data with higher hindsight TD error, better on-policiness and more accurate Q value should be assigned with higher weights during sampling. Thus most previous criteria only consider this strategy partially. We not only provide theoretical justifications for previous criteria, but also propose two new methods to compute the prioritization weight, namely ReMERN and ReMERT. ReMERN

learns an error network, while ReMERT exploits the temporal ordering of states. Both methods outperform previous prioritized sampling algorithms in challenging RL benchmarks, including MuJoCo, Atari and Meta-World.

[Relative Uncertainty Learning for Facial Expression Recognition](#)

- Yuhang Zhang, Chengrui Wang, Weihong Deng
- abstract@[open-review](#): In facial expression recognition (FER), the uncertainties introduced by inherent noises like ambiguous facial expressions and inconsistent labels raise concerns about the credibility of recognition results. To quantify these uncertainties and achieve good performance under noisy data, we regard uncertainty as a relative concept and propose an innovative uncertainty learning method called Relative Uncertainty Learning (RUL). Rather than assuming Gaussian uncertainty distributions for all datasets, RUL builds an extra branch to learn uncertainty from the relative difficulty of samples by feature mixup. Specifically, we use uncertainties as weights to mix facial features and design an add-up loss to encourage uncertainty learning. It is easy to implement and adds little or no extra computation overhead. Extensive experiments show that RUL outperforms state-of-the-art FER uncertainty learning methods in both real-world and synthetic noisy FER datasets. Besides, RUL also works well on other datasets such as CIFAR and Tiny ImageNet. The code is available at <https://github.com/zyh-uaiaaaa/Relative-Uncertainty-Learning>.

[An Information-theoretic Approach to Distribution Shifts](#)

- Marco Federici, Ryota Tomioka, Patrick Forré
- abstract@[open-review](#): Safely deploying machine learning models to the real world is often a challenging process. For example, models trained with data obtained from a specific geographic location tend to fail when queried with data obtained elsewhere, agents trained in a simulation can struggle to adapt when deployed in the real world or novel environments, and neural networks that are fit to a subset of the population might carry some selection bias into their decision process. In this work, we describe the problem of data shift from an information-theoretic perspective by (i) identifying and describing the different sources of error, (ii) comparing some of the most promising objectives explored in the recent domain generalization and fair classification literature. From our theoretical analysis and empirical evaluation, we conclude that the model selection procedure needs to be guided by careful considerations regarding the observed data, the factors used for correction, and the structure of the data-generating process.

[TRS: Transferability Reduced Ensemble via Promoting Gradient Diversity and Model Smoothness](#)

- Zhuolin Yang, Linyi Li, Xiaojun Xu, Shiliang Zuo, Qian Chen, Pan Zhou, Benjamin Rubinstein, Ce Zhang, Bo Li
- abstract@[open-review](#): Adversarial Transferability is an intriguing property - adversarial perturbation crafted against one model is also effective against another model, while these models are from different model families or training processes. To better protect ML systems against adversarial attacks, several questions are raised: what are the sufficient conditions for adversarial transferability, and how to bound it? Is there a way to reduce the adversarial transferability in order to improve the robustness of an ensemble ML model? To answer these questions, in this work we first theoretically analyze and outline sufficient conditions for adversarial transferability between models; then propose a practical algorithm to reduce the transferability between base models within an ensemble to improve its robustness. Our theoretical analysis shows that only promoting the orthogonality between gradients of base models is not enough to ensure low transferability; in the meantime, the model smoothness is an important factor to control the transferability. We also provide the lower and upper bounds of adversarial transferability under certain conditions. Inspired by our theoretical analysis, we propose an effective Transferability Reduced Smooth (TRS) ensemble training strategy to train a robust ensemble with low transferability by enforcing both gradient orthogonality and model smoothness between base models. We conduct extensive experiments on TRS and compare with 6 state-of-the-art ensemble baselines against 8 whitebox attacks on different datasets, demonstrating that the proposed TRS outperforms all baselines significantly.

[Towards Sample-Optimal Compressive Phase Retrieval with Sparse and Generative Priors](#)

- Zhaoqiang Liu, Subhroshekhar Ghosh, Jonathan Scarlett
- abstract@[open-review](#): Compressive phase retrieval is a popular variant of the standard compressive sensing problem in which the measurements only contain magnitude information. In this paper, motivated by recent advances in deep generative models, we provide recovery guarantees with near-optimal sample complexity for phase retrieval with generative priors. We first show that when using i.i.d. Gaussian measurements and an $\$L\$$ -Lipschitz continuous generative model with bounded $\$k\$$ -dimensional inputs, roughly $\$O(k \log L)$ samples suffice to guarantee that any signal minimizing an amplitude-based empirical loss function is close to the true signal. Attaining this sample complexity with a practical algorithm remains a difficult challenge, and finding a good initialization for gradient-based methods has been observed to pose a major bottleneck. To partially address this, we further show that roughly $\$O(k \log L)$ samples ensure sufficient closeness between the underlying signal and any $\{\text{em globally optimal}\}$ solution to an optimization problem designed for spectral initialization (though finding such a solution may still be challenging). We also adapt this result to sparse phase retrieval, and show that $\$O(s \log n)$ samples are sufficient for a similar guarantee when the underlying signal is $\$s\$$ -sparse and $\$n\$$ -dimensional, matching an information-theoretic lower bound. While these guarantees do not directly correspond to a practical algorithm, we propose a practical spectral initialization method motivated by our findings, and experimentally observe performance gains over various existing spectral initialization methods for sparse phase retrieval.

[Moser Flow: Divergence-based Generative Modeling on Manifolds](#)

- Noam Rozen, Aditya Grover, Maximilian Nickel, Yaron Lipman
- abstract@[open-review](#): We are interested in learning generative models for complex geometries described via manifolds, such as spheres, tori, and other implicit surfaces. Current extensions of existing (Euclidean) generative models are restricted to specific geometries and typically suffer from high computational costs. We introduce Moser Flow (MF), a new class of generative models within the family of continuous normalizing flows (CNF). MF also produces a CNF via a solution to the change-of-variable formula, however differently from other CNF methods, its model (learned) density is parameterized as the source (prior) density minus the divergence of a neural network (NN). The divergence is a local, linear differential operator, easy to approximate and calculate on manifolds. Therefore, unlike other CNFs, MF does not require invoking or backpropagating through an ODE solver during training. Furthermore, representing the model density explicitly as the divergence of a NN rather than as a solution of an ODE facilitates learning high fidelity densities. Theoretically, we prove that MF constitutes a universal density approximator under suitable assumptions. Empirically, we demonstrate for the first time the use of flow models for sampling from general curved surfaces and achieve significant improvements in density estimation, sample quality, and training complexity over existing CNFs on challenging synthetic geometries and real-world benchmarks from the earth and climate sciences.

[Structure-Aware Random Fourier Kernel for Graphs](#)

- Jinyuan Fang, Qiang Zhang, Zaiqiao Meng, Shangsong Liang
- abstract@[open-review](#): Gaussian Processes (GPs) define distributions over functions and their generalization capabilities depend heavily on the choice of kernels. In this paper, we propose a novel structure-aware random Fourier (SRF) kernel for GPs that brings several benefits when modeling graph-structured data. First, SRF kernel is defined with a spectral distribution based on the Fourier duality given by the Bochner's theorem, transforming the kernel learning problem to a distribution inference problem. Second, SRF kernel admits a random Fourier feature formulation that makes the kernel scalable for optimization. Third, SRF kernel enables to leverage geometric structures by taking subgraphs as inputs. To effectively optimize GPs with SRF kernel, we develop a variational EM algorithm, which alternates between an inference procedure (E-step) and a learning procedure (M-step). Experimental results on five real-world datasets show that our model can achieve state-of-the-art performance in two typical graph learning tasks, i.e., object classification and link prediction.

[Diffusion Schrödinger Bridge with Applications to Score-Based Generative Modeling](#)

- Valentin De Bortoli, James Thornton, Jeremy Heng, Arnaud Doucet
- abstract@[open-review](#): Progressively applying Gaussian noise transforms complex data distributions to approximately Gaussian. Reversing this dynamic defines a generative model. When the forward noising process is given by a Stochastic Differential Equation (SDE), Song et al (2021) demonstrate how the time inhomogeneous drift of the associated reverse-time SDE may be estimated using score-matching. A limitation of this approach is that the forward-time SDE must be run for a sufficiently long time for the final distribution to be approximately Gaussian. In contrast, solving the Schrödinger Bridge (SB) problem, i.e. an entropy-regularized optimal transport problem on path spaces, yields diffusions which generate samples from the data distribution in finite time. We present Diffusion SB (DSB), an original approximation of the Iterative Proportional Fitting (IPF) procedure to solve the SB problem, and provide theoretical analysis along with generative modeling experiments. The first DSB iteration recovers the methodology proposed by Song et al. (2021), with the flexibility of using shorter time intervals, as subsequent DSB iterations reduce the discrepancy between the final-time marginal of the forward (resp. backward) SDE with respect to the prior (resp. data) distribution. Beyond generative modeling, DSB offers a widely applicable computational optimal transport tool as the continuous state-space analogue of the popular Sinkhorn algorithm (Cuturi, 2013).

[Improving Transferability of Representations via Augmentation-Aware Self-Supervision](#)

- Hankook Lee, Kibok Lee, Kimin Lee, Honglak Lee, Jinwoo Shin
- abstract@[open-review](#): Recent unsupervised representation learning methods have shown to be effective in a range of vision tasks by learning representations invariant to data augmentations such as random cropping and color jittering. However, such invariance could be harmful to downstream tasks if they rely on the characteristics of the data augmentations, e.g., location- or color-sensitive. This is not an issue just for unsupervised learning; we found that this occurs even in supervised learning because it also learns to predict the same label for all augmented samples of an instance. To avoid such failures and obtain more generalizable representations, we suggest to optimize an auxiliary self-supervised loss, coined AugSelf, that learns the difference of augmentation parameters (e.g., cropping positions, color adjustment intensities) between two randomly augmented samples. Our intuition is that AugSelf encourages to preserve augmentation-aware information in learned representations, which could be beneficial for their transferability. Furthermore, AugSelf can easily be incorporated into recent state-of-the-art representation learning methods with a negligible additional training cost. Extensive experiments demonstrate that our simple idea consistently improves the transferability of representations learned by supervised and unsupervised methods in various transfer learning scenarios. The code is available at <https://github.com/hankook/AugSelf>.

[Long-Short Transformer: Efficient Transformers for Language and Vision](#)

- Chen Zhu, Wei Ping, Chaowei Xiao, Mohammad Shoeybi, Tom Goldstein, Anima Anandkumar, Bryan Catanzaro
- abstract@[open-review](#): Transformers have achieved success in both language and vision domains. However, it is prohibitively expensive to scale them to long sequences such as long documents or high-resolution images, because self-attention mechanism has quadratic time and memory complexities with respect to the input sequence length. In this paper, we propose Long-Short Transformer (Transformer-LS), an efficient self-attention mechanism for modeling long sequences with linear complexity for both language and vision tasks. It aggregates a novel long-range attention with dynamic projection to model distant correlations and a short-term attention to capture fine-grained local correlations. We propose a dual normalization strategy to account for the scale mismatch between the two attention mechanisms. Transformer-LS can be applied to both autoregressive and bidirectional models without additional complexity. Our method outperforms the state-of-the-art models on multiple tasks in language and vision domains, including the Long Range Arena benchmark, autoregressive language modeling, and ImageNet classification. For instance, Transformer-LS achieves 0.97 test BPC on enwik8 using half the number of parameters than previous method, while being faster and is able to handle 3x as long sequences compared to its full-attention version on the same hardware. On ImageNet, it can obtain the state-of-the-art results (e.g., a moderate size of 55.8M model solely trained on 224x224 ImageNet-1K can obtain Top-1 accuracy 84.1%), while being more scalable on high-resolution images. The source code and models are released at <https://github.com/NVIDIA/transformer-ls>.

[Post-Training Sparsity-Aware Quantization](#)

- Gil Shomron, Freddy Gabbay, Samer Kurzum, Uri Weiser
- abstract@[open-review](#): Quantization is a technique used in deep neural networks (DNNs) to increase execution performance and hardware efficiency. Uniform post-training quantization (PTQ) methods are common, since they can be implemented efficiently in hardware and do not require extensive hardware resources or a training set. Mapping FP32 models to INT8 using uniform PTQ yields models with negligible accuracy degradation; however, reducing precision below 8 bits with PTQ is challenging, as accuracy degradation becomes noticeable, due to the increase in quantization noise. In this paper, we propose a sparsity-aware quantization (SPARQ) method, in which the unstructured and dynamic activation sparsity is leveraged in different representation granularities. 4-bit quantization, for example, is employed by dynamically examining the bits of 8-bit values and choosing a window of 4 bits, while first skipping zero-value bits. Moreover, instead of quantizing activation-by-activation to 4 bits, we focus on pairs of 8-bit activations and examine whether one of the two is equal to zero. If one is equal to zero, the second can opportunistically use the other's 4-bit budget; if both do not equal zero, then each is dynamically quantized to 4 bits, as described. SPARQ achieves minor accuracy degradation and a practical hardware implementation.

[The Implicit Bias of Minima Stability: A View from Function Space](#)

- Rotem Mulayoff, Tomer Michaeli, Daniel Soudry
- abstract@[open-review](#): The loss terrains of over-parameterized neural networks have multiple global minima. However, it is well known that stochastic gradient descent (SGD) can stably converge only to minima that are sufficiently flat w.r.t. SGD's step size. In this paper we study the effect that this mechanism has on the function implemented by the trained model. First, we extend the existing knowledge on minima stability to non-differentiable minima, which are common in ReLU nets. We then use our stability results to study a single hidden layer univariate ReLU network. In this setting, we show that SGD is biased towards functions whose second derivative (w.r.t the input) has a bounded weighted $\|L\|_1$ norm, and this is regardless of the initialization. In particular, we show that the function implemented by the network upon convergence gets smoother as the learning rate increases. The weight multiplying the second derivative is larger around the center of the support of the training distribution, and smaller towards its boundaries, suggesting that a trained model tends to be smoother at the center of the training distribution.

[Breaking the Sample Complexity Barrier to Regret-Optimal Model-Free Reinforcement Learning](#)

- Gen Li, Laixi Shi, Yuxin Chen, Yuantao Gu, Yuejie Chi
- abstract@[open-review](#): Achieving sample efficiency in online episodic reinforcement learning (RL) requires optimally balancing exploration and exploitation. When it comes to a finite-horizon episodic Markov decision process with S states, A actions and horizon length H , substantial progress has been achieved towards characterizing the minimax-optimal regret, which scales on the order of $\sqrt{H^2SAT}$ (modulo log factors) with T the total number of samples. While several competing solution paradigms have been proposed to minimize regret, they are either memory-inefficient, or fall short of optimality unless the sample size exceeds an enormous threshold (e.g., $S^6A^4 \log(H)$ for existing model-free methods). To overcome such a large sample size barrier to efficient RL, we design a novel model-free algorithm, with space complexity $O(SAH)$, that achieves near-optimal regret as soon as the sample size exceeds the order of S^5A^3 . In terms of this sample size requirement (also referred to the initial burn-in cost), our method improves --- by at least a factor of S^5A^3 --- upon any prior memory-efficient algorithm that is asymptotically regret-optimal. Leveraging the recently introduced variance reduction strategy (also called ‘em reference-advantage decomposition’), the proposed algorithm employs an ‘em early-settled’ reference update rule, with the aid of two Q-learning sequences with upper and lower confidence bounds. The design

principle of our early-settled variance reduction method might be of independent interest to other RL settings that involve intricate exploration-exploitation trade-offs.

[Robust Auction Design in the Auto-bidding World](#)

- Santiago Balseiro, Yuan Deng, Jieming Mao, Vahab Mirrokni, Song Zuo
- abstract@[open-review](#): In classic auction theory, reserve prices are known to be effective for improving revenue for the auctioneer against quasi-linear utility maximizing bidders. The introduction of reserve prices, however, usually do not help improve total welfare of the auctioneer and the bidders. In this paper, we focus on value maximizing bidders with return on spend constraints---a paradigm that has drawn considerable attention recently as more advertisers adopt auto-bidding algorithms in advertising platforms---and show that the introduction of reserve prices has a novel impact on the market. Namely, by choosing reserve prices appropriately the auctioneer can improve not only the total revenue but also the total welfare. Our results also demonstrate that reserve prices are robust to bidder types, i.e., reserve prices work well for different bidder types, such as value maximizers and utility maximizers, without using bidder type information. We generalize these results for a variety of auction mechanisms such as VCG, GSP, and first-price auctions. Moreover, we show how to combine these results with additive boosts to improve the welfare of the outcomes of the auction further. Finally, we complement our theoretical observations with an empirical study confirming the effectiveness of these ideas using data from online advertising auctions.

[Weighted model estimation for offline model-based reinforcement learning](#)

- Toru Hishinuma, Kei Senda
- abstract@[open-review](#): This paper discusses model estimation in offline model-based reinforcement learning (MBRL), which is important for subsequent policy improvement using an estimated model. From the viewpoint of covariate shift, a natural idea is model estimation weighted by the ratio of the state-action distributions of offline data and real future data. However, estimating such a natural weight is one of the main challenges for off-policy evaluation, which is not easy to use. As an artificial alternative, this paper considers weighting with the state-action distribution ratio of offline data and simulated future data, which can be estimated relatively easily by standard density ratio estimation techniques for supervised learning. Based on the artificial weight, this paper defines a loss function for offline MBRL and presents an algorithm to optimize it. Weighting with the artificial weight is justified as evaluating an upper bound of the policy evaluation error. Numerical experiments demonstrate the effectiveness of weighting with the artificial weight.

[Practical, Provably-Correct Interactive Learning in the Realizable Setting: The Power of True Believers](#)

- Julian Katz-Samuels, Blake Mason, Kevin G. Jamieson, Rob Nowak
- abstract@[open-review](#): We consider interactive learning in the realizable setting and develop a general framework to handle problems ranging from best arm identification to active classification. We begin our investigation with the observation that agnostic algorithms \emph{cannot} be minimax-optimal in the realizable setting. Hence, we design novel computationally efficient algorithms for the realizable setting that match the minimax lower bound up to logarithmic factors and are general-purpose, accommodating a wide variety of function classes including kernel methods, Hölder smooth functions, and convex functions. The sample complexities of our algorithms can be quantified in terms of well-known quantities like the extended teaching dimension and haystack dimension. However, unlike algorithms based directly on those combinatorial quantities, our algorithms are computationally efficient. To achieve computational efficiency, our algorithms sample from the version space using Monte Carlo ``hit-and-run'' algorithms instead of maintaining the version space explicitly. Our approach has two key strengths. First, it is simple, consisting of two unifying, greedy algorithms. Second, our algorithms have the capability to seamlessly leverage prior knowledge that is often available and useful in practice. In addition to our new theoretical results, we demonstrate empirically that our algorithms are competitive with Gaussian process UCB methods.

[Deconditional Downscaling with Gaussian Processes](#)

- Siu Lun Chau, Shahine Bouabid, Dino Sejdinovic
- abstract@[open-review](#): Refining low-resolution (LR) spatial fields with high-resolution (HR) information, often known as statistical downscaling, is challenging as the diversity of spatial datasets often prevents direct matching of observations. Yet, when LR samples are modeled as aggregate conditional means of HR samples with respect to a mediating variable that is globally observed, the recovery of the underlying fine-grained field can be framed as taking an "inverse" of the conditional expectation, namely a deconditioning problem. In this work, we propose a Bayesian formulation of deconditioning which naturally recovers the initial reproducing kernel Hilbert space formulation from Hsu and Ramos (2019). We extend deconditioning to a downscaling setup and devise efficient conditional mean embedding estimator for multiresolution data. By treating conditional expectations as inter-domain features of the underlying field, a posterior for the latent field can be established as a solution to the deconditioning problem. Furthermore, we show that this solution can be viewed as a two-staged vector-valued kernel ridge regressor and show that it has a minimax optimal convergence rate under mild assumptions. Lastly, we demonstrate its proficiency in a synthetic and a real-world atmospheric field downscaling problem, showing substantial improvements over existing methods.

[Image Generation using Continuous Filter Atoms](#)

- Ze Wang, Seunghyun Hwang, Zichen Miao, Qiang Qiu
- abstract@[open-review](#): In this paper, we model the subspace of convolutional filters with a neural ordinary differential equation (ODE) to enable gradual changes in generated images. Decomposing convolutional filters over a set of filter atoms allows efficiently modeling and sampling from a subspace of high-dimensional filters. By further modeling filters atoms with a neural ODE, we show both empirically and theoretically that such introduced continuity can be propagated to the generated images, and thus achieves gradually evolved image generation. We support the proposed framework of image generation with continuous filter atoms using various experiments, including image-to-image translation and image generation conditioned on continuous labels. Without auxiliary network components and heavy supervision, the proposed continuous filter atoms allow us to easily manipulate the gradual change of generated images by controlling integration intervals of neural ordinary differential equation. This research sheds the light on using the subspace of network parameters to navigate the diverse appearance of image generation.

[Latent Equilibrium: A unified learning theory for arbitrarily fast computation with arbitrarily slow neurons](#)

- Paul Haider, Benjamin Ellenberger, Laura Kriener, Jakob Jordan, Walter Senn, Mihai A. Petrovici
- abstract@[open-review](#): The response time of physical computational elements is finite, and neurons are no exception. In hierarchical models of cortical networks each layer thus introduces a response lag. This inherent property of physical dynamical systems results in delayed processing of stimuli and causes a timing mismatch between network output and instructive signals, thus afflicting not only inference, but also learning. We introduce Latent Equilibrium, a new framework for inference and learning in networks of slow components which avoids these issues by harnessing the ability of biological neurons to phase-advance their output with respect to their membrane potential. This principle enables quasi-instantaneous inference independent of network depth and avoids the need for phased plasticity or computationally expensive network relaxation phases. We jointly derive disentangled neuron and synapse dynamics from a prospective energy function that depends on a network's generalized position and momentum. The resulting model can be interpreted as a biologically plausible approximation of error backpropagation in deep cortical networks with continuous-time, leaky neuronal dynamics and continuously active, local plasticity. We demonstrate successful learning of standard benchmark datasets, achieving competitive performance using both fully-connected and convolutional architectures, and show how our principle can be applied to detailed models of cortical microcircuitry. Furthermore, we study the robustness of our model to spatio-temporal substrate imperfections to demonstrate its feasibility for physical realization, be it *in vivo* or *in silico*.

Learning Fast-Inference Bayesian Networks

- Vaidyanathan Peruvemba Ramaswamy, Stefan Szeider
- abstract@[open-review](#): We propose new methods for learning Bayesian networks (BNs) that reliably support fast inference. We utilize maximum state space size as a more fine-grained measure for the BN's reasoning complexity than the standard treewidth measure, thereby accommodating the possibility that variables range over domains of different sizes. Our methods combine heuristic BN structure learning algorithms with the recently introduced MaxSAT-powered local improvement method (Peruvemba Ramaswamy and Szeider, AAAI'21). Our experiments show that our new learning methods produce BNs that support significantly faster exact probabilistic inference than BNs learned with treewidth bounds.

Per-Pixel Classification is Not All You Need for Semantic Segmentation

- Bowen Cheng, Alex Schwing, Alexander Kirillov
- abstract@[open-review](#): Modern approaches typically formulate semantic segmentation as a per-pixel classification task, while instance-level segmentation is handled with an alternative mask classification. Our key insight: mask classification is sufficiently general to solve both semantic- and instance-level segmentation tasks in a unified manner using the exact same model, loss, and training procedure. Following this observation, we propose MaskFormer, a simple mask classification model which predicts a set of binary masks, each associated with a single global class label prediction. Overall, the proposed mask classification-based method simplifies the landscape of effective approaches to semantic and panoptic segmentation tasks and shows excellent empirical results. In particular, we observe that MaskFormer outperforms per-pixel classification baselines when the number of classes is large. Our mask classification-based method outperforms both current state-of-the-art semantic (55.6 mIoU on ADE20K) and panoptic segmentation (52.7 PQ on COCO) models.

Deep Markov Factor Analysis: Towards Concurrent Temporal and Spatial Analysis of fMRI Data

- Amirreza Farnoosh, Sarah Ostadabbas
- abstract@[open-review](#): Factor analysis methods have been widely used in neuroimaging to transfer high dimensional imaging data into low dimensional, ideally interpretable representations. However, most of these methods overlook the highly nonlinear and complex temporal dynamics of neural processes when factorizing their imaging data. In this paper, we present deep Markov factor analysis (DMFA), a generative model that employs Markov property in a chain of low dimensional temporal embeddings together with spatial inductive assumptions, all related through neural networks, to capture temporal dynamics in functional magnetic resonance imaging (fMRI) data, and tackle their high spatial dimensionality, respectively. Augmented with a discrete latent, DMFA is able to cluster fMRI data in its low dimensional temporal embedding with regard to subject and cognitive state variability, therefore, enables validation of a variety of fMRI-driven neuroscientific hypotheses. Experimental results on both synthetic and real fMRI data demonstrate the capacity of DMFA in revealing interpretable clusters and capturing nonlinear temporal dependencies in these high dimensional imaging data.

BooVAE: Boosting Approach for Continual Learning of VAE

- Evgenii Egorov, Anna Kuzina, Evgeny Burnaev
- abstract@[open-review](#): Variational autoencoder (VAE) is a deep generative model for unsupervised learning, allowing to encode observations into the meaningful latent space. VAE is prone to catastrophic forgetting when tasks arrive sequentially, and only the data for the current one is available. We address this problem of continual learning for VAEs. It is known that the choice of the prior distribution over the latent space is crucial for VAE in the non-continual setting. We argue that it can also be helpful to avoid catastrophic forgetting. We learn the approximation of the aggregated posterior as a prior for each task. This approximation is parametrised as an additive mixture of distributions induced by an encoder evaluated at trainable pseudo-inputs. We use a greedy boosting-like approach with entropy regularisation to learn the components. This method encourages components diversity, which is essential as we aim at memorising the current task with the fewest components possible. Based on the learnable prior, we introduce an end-to-end approach for continual learning of VAEs and provide empirical studies on commonly used benchmarks (MNIST, Fashion MNIST, NotMNIST) and CelebA datasets. For each dataset, the proposed method avoids catastrophic forgetting in a fully automatic way.

Handling Long-tailed Feature Distribution in AdderNets

- Minjing Dong, Yunhe Wang, Xinghao Chen, Chang Xu
- abstract@[open-review](#): Adder neural networks (ANNs) are designed for low energy cost which replace expensive multiplications in convolutional neural networks (CNNs) with cheaper additions to yield energy-efficient neural networks and hardware accelerations. Although ANNs achieve satisfactory efficiency, there exist gaps between ANNs and CNNs where the accuracy of ANNs can hardly be compared to CNNs without the assistance of other training tricks, such as knowledge distillation. The inherent discrepancy lies in the similarity measurement between filters and features, however how to alleviate this difference remains unexplored. To locate the potential problem of ANNs, we focus on the property difference due to similarity measurement. We demonstrate that unordered heavy tails in ANNs could be the key component which prevents ANNs from achieving superior classification performance since fatter tails tend to overlap in feature space. Through pre-defining Multivariate Skew Laplace distributions and embedding feature distributions into the loss function, ANN features can be fully controlled and designed for various properties. We further present a novel method for tackling existing heavy tails in ANNs with only a modification of classifier where ANN features are clustered with their tails well-formulated through proposed angle-based constraint on the distribution parameters to encourage high diversity of tails. Experiments conducted on several benchmarks and comparison with other distributions demonstrate the effectiveness of proposed approach for boosting the performance of ANNs.

Pessimism Meets Invariance: Provably Efficient Offline Mean-Field Multi-Agent RL

- Minshuo Chen, Yan Li, Ethan Wang, Zhuoran Yang, Zhaoran Wang, Tuo Zhao
- abstract@[open-review](#): Mean-Field Multi-Agent Reinforcement Learning (MF-MARL) is attractive in the applications involving a large population of homogeneous agents, as it exploits the permutation invariance of agents and avoids the curse of many agents. Most existing results only focus on online settings, in which agents can interact with the environment during training. In some applications such as social welfare optimization, however, the interaction during training can be prohibitive or even unethical in the societal systems. To bridge such a gap, we propose a SAFARI (peSSimistic meAn-Field vAlue iteRatIon) algorithm for off-line MF-MARL, which only requires a handful of pre-collected experience data. Theoretically, under a weak coverage assumption that the experience dataset contains enough information about the optimal policy, we prove that for an episodic mean-field MDP with a horizon H and N training trajectories, SAFARI attains a sub-optimality gap of $\mathcal{O}(H^2 d_{\text{rm eff}} \sqrt{N})$, where $d_{\text{rm eff}}$ is the effective dimension of the function class for parameterizing the value function, but independent on the number of agents. Numerical experiments are provided.

A Law of Iterated Logarithm for Multi-Agent Reinforcement Learning

- Gugan Chandrashekhar Thoppe, Bhumesh Kumar
- abstract@[open-review](#): In Multi-Agent Reinforcement Learning (MARL), multiple agents interact with a common environment, as also with each other, for solving a shared problem in sequential decision-making. It has wide-ranging applications in gaming, robotics, finance, communication, etc. In this work, we derive a novel law of iterated logarithm for a family of distributed nonlinear stochastic approximation schemes that is useful in MARL. In particular, our result describes the convergence rate on almost every sample path where the algorithm converges. This result is the first of its kind in the

distributed setup and provides deeper insights than the existing ones, which only discuss convergence rates in the expected or the CLT sense. Importantly, our result holds under significantly weaker assumptions: neither the gossip matrix needs to be doubly stochastic nor the stepsizes square summable. As an application, we show that, for the stepsize $n^{-\gamma}$ with $\gamma \in (0, 1)$, the distributed TD(0) algorithm with linear function approximation has a convergence rate of $O(\sqrt{n^{-\gamma} \ln n})$ a.s.; for the $1/n$ type stepsize, the same is $O(\sqrt{n^{-1} \ln \ln n})$ a.s. These decay rates do not depend on the graph depicting the interactions among the different agents.

[MOMA: Multi-Object Multi-Actor Activity Parsing](#)

- Zelun Luo, Wanze Xie, Siddharth Kapoor, Yiyun Liang, Michael Cooper, Juan Carlos Niebles, Ehsan Adeli, Fei-Fei Li
- abstract@[open-review](#): Complex activities often involve multiple humans utilizing different objects to complete actions (e.g., in healthcare settings, physicians, nurses, and patients interact with each other and various medical devices). Recognizing activities poses a challenge that requires a detailed understanding of actors' roles, objects' affordances, and their associated relationships. Furthermore, these purposeful activities are composed of multiple achievable steps, including sub-activities and atomic actions, which jointly define a hierarchy of action parts. This paper introduces Activity Parsing as the overarching task of temporal segmentation and classification of activities, sub-activities, atomic actions, along with an instance-level understanding of actors, objects, and their relationships in videos. Involving multiple entities (actors and objects), we argue that traditional pair-wise relationships, often used in scene or action graphs, do not appropriately represent the dynamics between them. Hence, we introduce Action Hypergraph, a spatial-temporal graph containing hyperedges (i.e., edges with higher-order relationships), as a new representation. In addition, we introduce Multi-Object Multi-Actor (MOMA), the first benchmark and dataset dedicated to activity parsing. Lastly, to parse a video, we propose the HyperGraph Activity Parsing (HGAP) network, which outperforms several baselines, including those based on regular graphs and raw video data.

[The Pareto Frontier of model selection for general Contextual Bandits](#)

- Teodor Vanislavov Marinov, Julian Zimmert
- abstract@[open-review](#): Recent progress in model selection raises the question of the fundamental limits of these techniques. Under specific scrutiny has been model selection for general contextual bandits with nested policy classes, resulting in a COLT2020 open problem. It asks whether it is possible to obtain simultaneously the optimal single algorithm guarantees over all policies in a nested sequence of policy classes, or if otherwise this is possible for a trade-off $\alpha \in [1/2, 1]$ between complexity term and time: $\ln(\Pi_{ml})^{1-\alpha} T^\alpha$. We give a disappointing answer to this question. Even in the purely stochastic regime, the desired results are unobtainable. We present a Pareto frontier of up to logarithmic factors matching upper and lower bounds, thereby proving that an increase in the complexity term $\ln(\Pi_{ml})$ independent of T is unavoidable for general policy classes. As a side result, we also resolve a COLT2016 open problem concerning second-order bounds in full-information games.

[Teaching an Active Learner with Contrastive Examples](#)

- Chaoqi Wang, Adish Singla, Yuxin Chen
- abstract@[open-review](#): We study the problem of active learning with the added twist that the learner is assisted by a helpful teacher. We consider the following natural interaction protocol: At each round, the learner proposes a query asking for the label of an instance x^q , the teacher provides the requested label x^q, y^q along with explanatory information to guide the learning process. In this paper, we view this information in the form of an additional contrastive example (x^c, y^c) where x^c is picked from a set constrained by x^q (e.g., dissimilar instances with the same label). Our focus is to design a teaching algorithm that can provide an informative sequence of contrastive examples to the learner to speed up the learning process. We show that this leads to a challenging sequence optimization problem where the algorithm's choices at a given round depend on the history of interactions. We investigate an efficient teaching algorithm that adaptively picks these contrastive examples. We derive strong performance guarantees for our algorithm based on two problem-dependent parameters and further show that for specific types of active learners (e.g., a generalized binary search learner), the proposed teaching algorithm exhibits strong approximation guarantees. Finally, we illustrate our bounds and demonstrate the effectiveness of our teaching framework via two numerical case studies.

[Structured Denoising Diffusion Models in Discrete State-Spaces](#)

- Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, Rianne van den Berg
- abstract@[open-review](#): Denoising diffusion probabilistic models (DDPMs) [Ho et al. 2021] have shown impressive results on image and waveform generation in continuous state spaces. Here, we introduce Discrete Denoising Diffusion Probabilistic Models (D3PMs), diffusion-like generative models for discrete data that generalize the multinomial diffusion model of Hoogeboom et al. [2021], by going beyond corruption processes with uniform transition probabilities. This includes corruption with transition matrices that mimic Gaussian kernels in continuous space, matrices based on nearest neighbors in embedding space, and matrices that introduce absorbing states. The third allows us to draw a connection between diffusion models and autoregressive and mask-based generative models. We show that the choice of transition matrix is an important design decision that leads to improved results in image and text domains. We also introduce a new loss function that combines the variational lower bound with an auxiliary cross entropy loss. For text, this model class achieves strong results on character-level text generation while scaling to large vocabularies on LM1B. On the image dataset CIFAR-10, our models approach the sample quality and exceed the log-likelihood of the continuous-space DDPM model.

[Emergent Communication of Generalizations](#)

- Jesse Mu, Noah Goodman
- abstract@[open-review](#): To build agents that can collaborate effectively with others, recent research has trained artificial agents to communicate with each other in Lewis-style referential games. However, this often leads to successful but uninterpretable communication. We argue that this is due to the game objective: communicating about a single object in a shared visual context is prone to overfitting and does not encourage language useful beyond concrete reference. In contrast, human language conveys a rich variety of abstract ideas. To promote such skills, we propose games that require communicating generalizations over sets of objects representing abstract visual concepts, optionally with separate contexts for each agent. We find that these games greatly improve systematicity and interpretability of the learned languages, according to several metrics in the literature. Finally, we propose a method for identifying logical operations embedded in the emergent languages by learning an approximate compositional reconstruction of the language.

[Distributed Machine Learning with Sparse Heterogeneous Data](#)

- Dominic Richards, Sahand Negahban, Patrick Rebeschini
- abstract@[open-review](#): Motivated by distributed machine learning settings such as Federated Learning, we consider the problem of fitting a statistical model across a distributed collection of heterogeneous data sets whose similarity structure is encoded by a graph topology. Precisely, we analyse the case where each node is associated with fitting a sparse linear model, and edges join two nodes if the difference of their solutions is also sparse. We propose a method based on Basis Pursuit Denoising with a total variation penalty, and provide finite sample guarantees for sub-Gaussian design matrices. Taking the root of the tree as a reference node, we show that if the sparsity of the differences across nodes is smaller than the sparsity at the root, then recovery is successful with fewer samples than by solving the problems independently, or by using methods that rely on a large overlap in the signal supports, such as the group Lasso. We consider both the noiseless and noisy setting, and numerically investigate the performance of distributed methods based on Distributed Alternating Direction Methods of Multipliers (ADMM) and hyperspectral unmixing.

Manipulating SGD with Data Ordering Attacks

- I Shumailov, Zakhar Shumaylov, Dmitry Kazhdan, Yiren Zhao, Nicolas Papernot, Murat A. Erdogdu, Ross J Anderson
- abstract@[open-review](#): Machine learning is vulnerable to a wide variety of attacks. It is now well understood that by changing the underlying data distribution, an adversary can poison the model trained with it or introduce backdoors. In this paper we present a novel class of training-time attacks that require no changes to the underlying dataset or model architecture, but instead only change the order in which data are supplied to the model. In particular, we find that the attacker can either prevent the model from learning, or poison it to learn behaviours specified by the attacker. Furthermore, we find that even a single adversarially-ordered epoch can be enough to slow down model learning, or even to reset all of the learning progress. Indeed, the attacks presented here are not specific to the model or dataset, but rather target the stochastic nature of modern learning procedures. We extensively evaluate our attacks on computer vision and natural language benchmarks to find that the adversary can disrupt model training and even introduce backdoors.

Graph Posterior Network: Bayesian Predictive Uncertainty for Node Classification

- Maximilian Stadler, Bertrand Charpentier, Simon Geisler, Daniel Zügner, Stephan Günnemann
- abstract@[open-review](#): The interdependence between nodes in graphs is key to improve class prediction on nodes, utilized in approaches like Label Propagation (LP) or in Graph Neural Networks (GNNs). Nonetheless, uncertainty estimation for non-independent node-level predictions is under-explored. In this work, we explore uncertainty quantification for node classification in three ways: (1) We derive three axioms explicitly characterizing the expected predictive uncertainty behavior in homophilic attributed graphs.(2) We propose a new model Graph Posterior Network (GPN) which explicitly performs Bayesian posterior updates for predictions on interdependent nodes. GPN provably obeys the proposed axioms. (3) We extensively evaluate GPN and a strong set of baselines on semi-supervised node classification including detection of anomalous features, and detection of left-out classes. GPN outperforms existing approaches for uncertainty estimation in the experiments.

Locality Sensitive Teaching

- Zhaozhuo Xu, Beidi Chen, Chaojian Li, Weiyang Liu, Le Song, Yingyan Lin, Anshumali Shrivastava
- abstract@[open-review](#): The emergence of the Internet-of-Things (IoT) sheds light on applying the machine teaching (MT) algorithms for online personalized education on home devices. This direction becomes more promising during the COVID-19 pandemic when in-person education becomes infeasible. However, as one of the most influential and practical MT paradigms, iterative machine teaching (IMT) is prohibited on IoT devices due to its inefficient and unscalable algorithms. IMT is a paradigm where a teacher feeds examples iteratively and intelligently based on the learner's status. In each iteration, current IMT algorithms greedily traverse the whole training set to find an example for the learner, which is computationally expensive in practice. We propose a novel teaching framework, Locality Sensitive Teaching (LST), based on locality sensitive sampling, to overcome these challenges. LST has provable near-constant time complexity, which is exponentially better than the existing baseline. With at most 425.12x speedups and 99.76% energy savings over IMT, LST is the first algorithm that enables energy and time efficient machine teaching on IoT devices. Owing to LST's substantial efficiency and scalability, it is readily applicable in real-world education scenarios.

No-Press Diplomacy from Scratch

- Anton Bakhtin, David Wu, Adam Lerer, Noam Brown
- abstract@[open-review](#): Prior AI successes in complex games have largely focused on settings with at most hundreds of actions at each decision point. In contrast, Diplomacy is a game with more than 10^{20} possible actions per turn. Previous attempts to address games with large branching factors, such as Diplomacy, StarCraft, and Dota, used human data to bootstrap the policy or used handcrafted reward shaping. In this paper, we describe an algorithm for action exploration and equilibrium approximation in games with combinatorial action spaces. This algorithm simultaneously performs value iteration while learning a policy proposal network. A double oracle step is used to explore additional actions to add to the policy proposals. At each state, the target state value and policy for the model training are computed via an equilibrium search procedure. Using this algorithm, we train an agent, DORA, completely from scratch for a popular two-player variant of Diplomacy and show that it achieves superhuman performance. Additionally, we extend our methods to full-scale no-press Diplomacy and for the first time train an agent from scratch with no human data. We present evidence that this agent plays a strategy that is incompatible with human-data bootstrapped agents. This presents the first strong evidence of multiple equilibria in Diplomacy and suggests that self play alone may be insufficient for achieving superhuman performance in Diplomacy.

Remember What You Want to Forget: Algorithms for Machine Unlearning

- Ayush Sekhari, Jayadev Acharya, Gautam Kamath, Ananda Theertha Suresh
- abstract@[open-review](#): We study the problem of unlearning datapoints from a learnt model. The learner first receives a dataset \mathcal{S} drawn i.i.d. from an unknown distribution, and outputs a model w that performs well on unseen samples from the same distribution. However, at some point in the future, any training datapoint $z \in \mathcal{S}$ can request to be unlearned, thus prompting the learner to modify its output model while still ensuring the same accuracy guarantees. We initiate a rigorous study of generalization in machine unlearning, where the goal is to perform well on previously unseen datapoints. Our focus is on both computational and storage complexity. For the setting of convex losses, we provide an unlearning algorithm that can unlearn up to $O(n/d^{1/4})$ samples, where d is the problem dimension. In comparison, in general, differentially private learning (which implies unlearning) only guarantees deletion of $O(n/d^{1/2})$ samples. This demonstrates a novel separation between differential privacy and machine unlearning.

Learning latent causal graphs via mixture oracles

- Bohdan Kivva, Goutham Rajendran, Pradeep Ravikumar, Bryon Aragam
- abstract@[open-review](#): We study the problem of reconstructing a causal graphical model from data in the presence of latent variables. The main problem of interest is recovering the causal structure over the latent variables while allowing for general, potentially nonlinear dependencies. In many practical problems, the dependence between raw observations (e.g. pixels in an image) is much less relevant than the dependence between certain high-level, latent features (e.g. concepts or objects), and this is the setting of interest. We provide conditions under which both the latent representations and the underlying latent causal model are identifiable by a reduction to a mixture oracle. These results highlight an intriguing connection between the well-studied problem of learning the order of a mixture model and the problem of learning the bipartite structure between observables and unobservables. The proof is constructive, and leads to several algorithms for explicitly reconstructing the full graphical model. We discuss efficient algorithms and provide experiments illustrating the algorithms in practice.

ErrorCompensatedX: error compensation for variance reduced algorithms

- Hanlin Tang, Yao Li, Ji Liu, Ming Yan
- abstract@[open-review](#): Communication cost is one major bottleneck for the scalability for distributed learning. One approach to reduce the communication cost is to compress the gradient during communication. However, directly compressing the gradient decelerates the convergence speed, and the resulting algorithm may diverge for biased compression. Recent work addressed this problem for stochastic gradient descent by adding back the compression error from the previous step. This idea was further extended to one class of variance reduced algorithms, where the variance of the stochastic gradient is reduced by taking a moving average over all history gradients. However, our analysis shows that just adding the previous step's compression

error, as done in existing work, does not fully compensate the compression error. So, we propose ErrorCompensateX, which uses the compression error from the previous two steps. We show that ErrorCompensateX can achieve the same asymptotic convergence rate with the training without compression. Moreover, we provide a unified theoretical analysis framework for this class of variance reduced algorithms, with or without error compensation.

[Deep Contextual Video Compression](#)

- Jiahao Li, Bin Li, Yan Lu
- abstract@[open-review](#): Most of the existing neural video compression methods adopt the predictive coding framework, which first generates the predicted frame and then encodes its residue with the current frame. However, as for compression ratio, predictive coding is only a sub-optimal solution as it uses simple subtraction operation to remove the redundancy across frames. In this paper, we propose a deep contextual video compression framework to enable a paradigm shift from predictive coding to conditional coding. In particular, we try to answer the following questions: how to define, use, and learn condition under a deep video compression framework. To tap the potential of conditional coding, we propose using feature domain context as condition. This enables us to leverage the high dimension context to carry rich information to both the encoder and the decoder, which helps reconstruct the high-frequency contents for higher video quality. Our framework is also extensible, in which the condition can be flexibly designed. Experiments show that our method can significantly outperform the previous state-of-the-art (SOTA) deep video compression methods. When compared with x265 using veryslow preset, we can achieve 26.0% bitrate saving for 1080P standard test videos.

[On the Frequency Bias of Generative Models](#)

- Katja Schwarz, Yiyi Liao, Andreas Geiger
- abstract@[open-review](#): The key objective of Generative Adversarial Networks (GANs) is to generate new data with the same statistics as the provided training data. However, multiple recent works show that state-of-the-art architectures yet struggle to achieve this goal. In particular, they report an elevated amount of high frequencies in the spectral statistics which makes it straightforward to distinguish real and generated images. Explanations for this phenomenon are controversial: While most works attribute the artifacts to the generator, other works point to the discriminator. We take a sober look at those explanations and provide insights on what makes proposed measures against high-frequency artifacts effective. To achieve this, we first independently assess the architectures of both the generator and discriminator and investigate if they exhibit a frequency bias that makes learning the distribution of high-frequency content particularly problematic. Based on these experiments, we make the following four observations: 1) Different upsampling operations bias the generator towards different spectral properties. 2) Checkerboard artifacts introduced by upsampling cannot explain the spectral discrepancies alone as the generator is able to compensate for these artifacts. 3) The discriminator does not struggle with detecting high frequencies per se but rather struggles with frequencies of low magnitude. 4) The downsampling operations in the discriminator can impair the quality of the training signal it provides. In light of these findings, we analyze proposed measures against high-frequency artifacts in state-of-the-art GAN training but find that none of the existing approaches can fully resolve spectral artifacts yet. Our results suggest that there is great potential in improving the discriminator and that this could be key to match the distribution of the training data more closely.

[Learning curves of generic features maps for realistic datasets with a teacher-student model](#)

- Bruno Loureiro, Cedric Gerbelot, Hugo Cui, Sebastian Goldt, Florent Krzakala, Marc Mezard, Lenka Zdeborová
- abstract@[open-review](#): Teacher-student models provide a framework in which the typical-case performance of high-dimensional supervised learning can be described in closed form. The assumptions of Gaussian i.i.d. input data underlying the canonical teacher-student model may, however, be perceived as too restrictive to capture the behaviour of realistic data sets. In this paper, we introduce a Gaussian covariate generalisation of the model where the teacher and student can act on different spaces, generated with fixed, but generic feature maps. While still solvable in a closed form, this generalization is able to capture the learning curves for a broad range of realistic data sets, thus redeeming the potential of the teacher-student framework. Our contribution is then two-fold: first, we prove a rigorous formula for the asymptotic training loss and generalisation error. Second, we present a number of situations where the learning curve of the model captures the one of a realistic data set learned with kernel regression and classification, with out-of-the-box feature maps such as random projections or scattering transforms, or with pre-learned ones - such as the features learned by training multi-layer neural networks. We discuss both the power and the limitations of the framework.

[It Has Potential: Gradient-Driven Denoisers for Convergent Solutions to Inverse Problems](#)

- Regev Cohen, Yochai Blau, Daniel Freedman, Ehud Rivlin
- abstract@[open-review](#): In recent years there has been increasing interest in leveraging denoisers for solving general inverse problems. Two leading frameworks are regularization-by-denoising (RED) and plug-and-play priors (PnP) which incorporate explicit likelihood functions with priors induced by denoising algorithms. RED and PnP have shown state-of-the-art performance in diverse imaging tasks when powerful denoisers are used, such as convolutional neural networks (CNNs). However, the study of their convergence remains an active line of research. Recent works derive the convergence of RED and PnP methods by treating CNN denoisers as approximations for maximum a posteriori (MAP) or minimum mean square error (MMSE) estimators. Yet, state-of-the-art denoisers cannot be interpreted as either MAP or MMSE estimators, since they typically do not exhibit symmetric Jacobians. Furthermore, obtaining stable inverse algorithms often requires controlling the Lipschitz constant of CNN denoisers during training. Precisely enforcing this constraint is impractical, hence, convergence cannot be completely guaranteed. In this work, we introduce image denoisers derived as the gradients of smooth scalar-valued deep neural networks, acting as potentials. This ensures two things: (1) the proposed denoisers display symmetric Jacobians, allowing for MAP and MMSE estimators interpretation; (2) the denoisers may be integrated into RED and PnP schemes with backtracking step size, removing the need for enforcing their Lipschitz constant. To show the latter, we develop a simple inversion method that utilizes the proposed denoisers. We theoretically establish its convergence to stationary points of an underlying objective function consisting of the learned potentials. We numerically validate our method through various imaging experiments, showing improved results compared to standard RED and PnP methods, and with additional provable stability.

[Training Over-parameterized Models with Non-decomposable Objectives](#)

- Harikrishna Narasimhan, Aditya K. Menon
- abstract@[open-review](#): Many modern machine learning applications come with complex and nuanced design goals such as minimizing the worst-case error, satisfying a given precision or recall target, or enforcing group-fairness constraints. Popular techniques for optimizing such non-decomposable objectives reduce the problem into a sequence of cost-sensitive learning tasks, each of which is then solved by re-weighting the training loss with example-specific costs. We point out that the standard approach of re-weighting the loss to incorporate label costs can produce unsatisfactory results when used to train over-parameterized models. As a remedy, we propose new cost-sensitive losses that extend the classical idea of logit adjustment to handle more general cost matrices. Our losses are calibrated, and can be further improved with distilled labels from a teacher model. Through experiments on benchmark image datasets, we showcase the effectiveness of our approach in training ResNet models with common robust and constrained optimization objectives.

[Reinforcement learning for optimization of variational quantum circuit architectures](#)

- Mateusz Ostaszewski, Lea M. Trenkwalder, Wojciech Masarczyk, Eleanor Scerri, Vedran Dunjko
- abstract@[open-review](#): The study of Variational Quantum Eigensolvers (VQEs) has been in the spotlight in recent times as they may lead to real-world applications of near-term quantum devices. However, their performance depends on the structure of the used variational ansatz, which requires balancing

the depth and expressivity of the corresponding circuit. At the same time, near-term restrictions limit the depth of the circuit we can expect to run. Thus, the optimization of the VQE ansatz requires maximizing the expressivity of the circuit while maintaining low depth. In recent years, various methods for VQE structure optimization have been introduced but the capacities of machine learning to aid with this problem have not yet been extensively investigated. In this work, we propose a reinforcement learning algorithm that autonomously explores the space of possible ansatze, identifying economic circuits which still yield accurate ground energy estimates. The algorithm uses a feedback-driven curriculum learning method that autonomously adapts the complexity of the learning problem to the current performance of the learning algorithm and it incrementally improves the accuracy of the result while minimizing the circuit depth. We showcase the performance of our algorithm on the problem of estimating the ground-state energy of lithium hydride (LiH) in various configurations. In this well-known benchmark problem, we achieve chemical accuracy and state-of-the-art results in terms of circuit depth.

[Moshpit SGD: Communication-Efficient Decentralized Training on Heterogeneous Unreliable Devices](#)

- Max Ryabinin, Eduard Gorbunov, Vsevolod Plokhotnyuk, Gennady Pekhimenko
- abstract@[open-review](#): Training deep neural networks on large datasets can often be accelerated by using multiple compute nodes. This approach, known as distributed training, can utilize hundreds of computers via specialized message-passing protocols such as Ring All-Reduce. However, running these protocols at scale requires reliable high-speed networking that is only available in dedicated clusters. In contrast, many real-world applications, such as federated learning and cloud-based distributed training, operate on unreliable devices with unstable network bandwidth. As a result, these applications are restricted to using parameter servers or gossip-based averaging protocols. In this work, we lift that restriction by proposing Moshpit All-Reduce — an iterative averaging protocol that exponentially converges to the global average. We demonstrate the efficiency of our protocol for distributed optimization with strong theoretical guarantees. The experiments show 1.3x speedup for ResNet-50 training on ImageNet compared to competitive gossip-based strategies and 1.5x speedup when training ALBERT-large on preemptible compute nodes.

[IRM—when it works and when it doesn't: A test case of natural language inference](#)

- Yana Dranker, He He, Yonatan Belinkov
- abstract@[open-review](#): Invariant Risk Minimization (IRM) is a recently proposed framework for out-of-distribution (o.o.d) generalization. Most of the studies on IRM so far have focused on theoretical results, toy problems, and simple models. In this work, we investigate the applicability of IRM to bias mitigation—a special case of o.o.d generalization—in increasingly naturalistic settings and deep models. Using natural language inference (NLI) as a test case, we start with a setting where both the dataset and the bias are synthetic, continue with a natural dataset and synthetic bias, and end with a fully realistic setting with natural datasets and bias. Our results show that in naturalistic settings, learning complex features in place of the bias proves to be difficult, leading to a rather small improvement over empirical risk minimization. Moreover, we find that in addition to being sensitive to random seeds, the performance of IRM also depends on several critical factors, notably dataset size, bias prevalence, and bias strength, thus limiting IRM's advantage in practical scenarios. Our results highlight key challenges in applying IRM to real-world scenarios, calling for a more naturalistic characterization of the problem setup for o.o.d generalization.

[Self-Supervised Learning Disentangled Group Representation as Feature](#)

- Tan Wang, Zhongqi Yue, Jianqiang Huang, Qianru Sun, Hanwang Zhang
- abstract@[open-review](#): A good visual representation is an inference map from observations (images) to features (vectors) that faithfully reflects the hidden modularized generative factors (semantics). In this paper, we formulate the notion of "good" representation from a group-theoretic view using Higgins' definition of disentangled representation, and show that existing Self-Supervised Learning (SSL) only disentangles simple augmentation features such as rotation and colorization, thus unable to modularize the remaining semantics. To break the limitation, we propose an iterative SSL algorithm: Iterative Partition-based Invariant Risk Minimization (IP-IRM), which successfully grounds the abstract semantics and the group acting on them into concrete contrastive learning. At each iteration, IP-IRM first partitions the training samples into two subsets that correspond to an entangled group element. Then, it minimizes a subset-invariant contrastive loss, where the invariance guarantees to disentangle the group element. We prove that IP-IRM converges to a fully disentangled representation and show its effectiveness on various benchmarks. Codes are available at <https://github.com/Wang-CN/IP-IRM>.

[SalKG: Learning From Knowledge Graph Explanations for Commonsense Reasoning](#)

- Aaron Chan, Jiashu Xu, Boyuan Long, Soumya Sanyal, Tanishq Gupta, Xiang Ren
- abstract@[open-review](#): Augmenting pre-trained language models with knowledge graphs (KGs) has achieved success on various commonsense reasoning tasks. However, for a given task instance, the KG, or certain parts of the KG, may not be useful. Although KG-augmented models often use attention to focus on specific KG components, the KG is still always used, and the attention mechanism is never explicitly taught which KG components should be used. Meanwhile, saliency methods can measure how much a KG feature (e.g., graph, node, path) influences the model to make the correct prediction, thus explaining which KG features are useful. This paper explores how saliency explanations can be used to improve KG-augmented models' performance. First, we propose to create coarse (Is the KG useful?) and fine (Which nodes/paths in the KG are useful?) saliency explanations. Second, to motivate saliency-based supervision, we analyze oracle KG-augmented models which directly use saliency explanations as extra inputs for guiding their attention. Third, we propose SalKG, a framework for KG-augmented models to learn from coarse and/or fine saliency explanations. Given saliency explanations created from a task's training set, SalKG jointly trains the model to predict the explanations, then solve the task by attending to KG features highlighted by the predicted explanations. On three popular commonsense QA benchmarks (CSQA, OBQA, CODAH) and a range of KG-augmented models, we show that SalKG can yield considerable performance gains --- up to 2.76% absolute improvement on CSQA.

[Supervising the Transfer of Reasoning Patterns in VQA](#)

- Corentin Kervadec, Christian Wolf, Grigory Antipov, Moez Baccouche, Madiha Nadri
- abstract@[open-review](#): Methods for Visual Question Answering (VQA) are notorious for leveraging dataset biases rather than performing reasoning, hindering generalization. It has been recently shown that better reasoning patterns emerge in attention layers of a state-of-the-art VQA model when they are trained on perfect (oracle) visual inputs. This provides evidence that deep neural networks can learn to reason when training conditions are favorable enough. However, transferring this learned knowledge to deployable models is a challenge, as much of it is lost during the transfer. We propose a method for knowledge transfer based on a regularization term in our loss function, supervising the sequence of required reasoning operations. We provide a theoretical analysis based on PAC-learning, showing that such program prediction can lead to decreased sample complexity under mild hypotheses. We also demonstrate the effectiveness of this approach experimentally on the GQA dataset and show its complementarity to BERT-like self-supervised pre-training.

[Conformal Bayesian Computation](#)

- Edwin Fong, Chris C Holmes
- abstract@[open-review](#): We develop scalable methods for producing conformal Bayesian predictive intervals with finite sample calibration guarantees. Bayesian posterior predictive distributions, $p(y \mid x)$, characterize subjective beliefs on outcomes of interest, y , conditional on predictors, x . Bayesian prediction is well-calibrated when the model is true, but the predictive intervals may exhibit poor empirical coverage when the model is misspecified, under the so called \mathcal{M} -open perspective. In contrast, conformal inference provides finite sample frequentist guarantees on predictive confidence intervals without the requirement of model fidelity. Using 'add-one-in' importance sampling, we show that conformal Bayesian

predictive intervals are efficiently obtained from re-weighted posterior samples of model parameters. Our approach contrasts with existing conformal methods that require expensive refitting of models or data-splitting to achieve computational efficiency. We demonstrate the utility on a range of examples including extensions to partially exchangeable settings such as hierarchical models.

[A Unified Approach to Fair Online Learning via Blackwell Approachability](#)

- Evgenii Chzhen, Christophe Giraud, Gilles Stoltz
- abstract@[open-review](#): We provide a setting and a general approach to fair online learning with stochastic sensitive and non-sensitive contexts. The setting is a repeated game between the Player and Nature, where at each stage both pick actions based on the contexts. Inspired by the notion of unawareness, we assume that the Player can only access the non-sensitive context before making a decision, while we discuss both cases of Nature accessing the sensitive contexts and Nature unaware of the sensitive contexts. Adapting Blackwell's approachability theory to handle the case of an unknown contexts' distribution, we provide a general necessary and sufficient condition for learning objectives to be compatible with some fairness constraints. This condition is instantiated on (group-wise) no-regret and (group-wise) calibration objectives, and on demographic parity as an additional constraint. When the objective is not compatible with the constraint, the provided framework permits to characterise the optimal trade-off between the two.

[Training Neural Networks is ER-complete](#)

- Mikkel Abrahamsen, Linda Kleist, Tillmann Miltzow
- abstract@[open-review](#): Given a neural network, training data, and a threshold, finding weights for the neural network such that the total error is below the threshold is known to be NP-hard. We determine the algorithmic complexity of this fundamental problem precisely, by showing that it is $\exists \mathbb{R} \text{-complete}$. This means that the problem is equivalent, up to polynomial time reductions, to deciding whether a system of polynomial equations and inequalities with integer coefficients and real unknowns has a solution. If, as widely expected, $\exists \mathbb{R}$ is strictly larger than NP, our work implies that the problem of training neural networks is not even in NP. Neural networks are usually trained using some variation of backpropagation. The result of this paper gives an explanation why techniques commonly used to solve big instances of NP-complete problems (such as SAT solvers, IP solvers, local search, dynamic programming, etc.) seem to be of no use to this task.

[Understanding the Under-Coverage Bias in Uncertainty Estimation](#)

- Yu Bai, Song Mei, Huan Wang, Caiming Xiong
- abstract@[open-review](#): Estimating the data uncertainty in regression tasks is often done by learning a quantile function or a prediction interval of the true label conditioned on the input. It is frequently observed that quantile regression---a vanilla algorithm for learning quantiles with asymptotic guarantees---tends to *under-cover* than the desired coverage level in reality. While various fixes have been proposed, a more fundamental understanding of why this under-coverage bias happens in the first place remains elusive. In this paper, we present a rigorous theoretical study on the coverage of uncertainty estimation algorithms in learning quantiles. We prove that quantile regression suffers from an inherent under-coverage bias, in a vanilla setting where we learn a realizable linear quantile function and there is more data than parameters. More quantitatively, for $\alpha > 0.5$ and small d/n , the α -quantile learned by quantile regression roughly achieves coverage $\alpha - (\alpha - 1/2) \cdot d/n$ regardless of the noise distribution, where d is the input dimension and n is the number of training data. Our theory reveals that this under-coverage bias stems from a certain high-dimensional parameter estimation error that is not implied by existing theories on quantile regression. Experiments on simulated and real data verify our theory and further illustrate the effect of various factors such as sample size and model capacity on the under-coverage bias in more practical setups.

[Decentralized Q-learning in Zero-sum Markov Games](#)

- Muhammed Sayin, Kaiqing Zhang, David Leslie, Tamer Basar, Asuman Ozdaglar
- abstract@[open-review](#): We study multi-agent reinforcement learning (MARL) in infinite-horizon discounted zero-sum Markov games. We focus on the practical but challenging setting of decentralized MARL, where agents make decisions without coordination by a centralized controller, but only based on their own payoffs and local actions executed. The agents need not observe the opponent's actions or payoffs, possibly being even oblivious to the presence of the opponent, nor be aware of the zero-sum structure of the underlying game, a setting also referred to as radically uncoupled in the literature of learning in games. In this paper, we develop a radically uncoupled Q-learning dynamics that is both rational and convergent: the learning dynamics converges to the best response to the opponent's strategy when the opponent follows an asymptotically stationary strategy; when both agents adopt the learning dynamics, they converge to the Nash equilibrium of the game. The key challenge in this decentralized setting is the non-stationarity of the environment from an agent's perspective, since both her own payoffs and the system evolution depend on the actions of other agents, and each agent adapts her policies simultaneously and independently. To address this issue, we develop a two-timescale learning dynamics where each agent updates her local Q-function and value function estimates concurrently, with the latter happening at a slower timescale.

[Fast Certified Robust Training with Short Warmup](#)

- Zhouxing Shi, Yihan Wang, Huan Zhang, Jinfeng Yi, Cho-Jui Hsieh
- abstract@[open-review](#): Recently, bound propagation based certified robust training methods have been proposed for training neural networks with certifiable robustness guarantees. Despite that state-of-the-art (SOTA) methods including interval bound propagation (IBP) and CROWN-IBP have per-batch training complexity similar to standard neural network training, they usually use a long warmup schedule with hundreds or thousands epochs to reach SOTA performance and are thus still costly. In this paper, we identify two important issues in existing methods, namely exploded bounds at initialization, and the imbalance in ReLU activation states and improve IBP training. These two issues make certified training difficult and unstable, and thereby long warmup schedules were needed in prior works. To mitigate these issues and conduct faster certified training with shorter warmup, we propose three improvements based on IBP training: 1) We derive a new weight initialization method for IBP training; 2) We propose to fully add Batch Normalization (BN) to each layer in the model, since we find BN can reduce the imbalance in ReLU activation states; 3) We also design regularization to explicitly tighten certified bounds and balance ReLU activation states during wamrup. We are able to obtain 65.03% verified error on CIFAR-10 ($\epsilon = \frac{8}{255}$) and 82.36% verified error on TinyImageNet ($\epsilon = \frac{1}{255}$) using very short training schedules (160 and 80 total epochs, respectively), outperforming literature SOTA trained with hundreds or thousands epochs under the same network architecture. The code is available at <https://github.com/shizhouxing/Fast-Certified-Robust-Training>.

[Vector-valued Distance and Gyrocalculus on the Space of Symmetric Positive Definite Matrices](#)

- Federico López, Beatrice Pozzetti, Steve Trettel, Michael Strube, Anna Wienhard
- abstract@[open-review](#): We propose the use of the vector-valued distance to compute distances and extract geometric information from the manifold of symmetric positive definite matrices (SPD), and develop gyrovector calculus, constructing analogs of vector space operations in this curved space. We implement these operations and showcase their versatility in the tasks of knowledge graph completion, item recommendation, and question answering. In experiments, the SPD models outperform their equivalents in Euclidean and hyperbolic space. The vector-valued distance allows us to visualize embeddings, showing that the models learn to disentangle representations of positive samples from negative ones.

[Improved Transformer for High-Resolution GANs](#)

- Long Zhao, Zizhao Zhang, Ting Chen, Dimitris Metaxas, Han Zhang
- abstract@[open-review](#): Attention-based models, exemplified by the Transformer, can effectively model long range dependency, but suffer from the quadratic complexity of self-attention operation, making them difficult to be adopted for high-resolution image generation based on Generative Adversarial Networks (GANs). In this paper, we introduce two key ingredients to Transformer to address this challenge. First, in low-resolution stages of the generative process, standard global self-attention is replaced with the proposed multi-axis blocked self-attention which allows efficient mixing of local and global attention. Second, in high-resolution stages, we drop self-attention while only keeping multi-layer perceptrons reminiscent of the implicit neural function. To further improve the performance, we introduce an additional self-modulation component based on cross-attention. The resulting model, denoted as HiT, has a nearly linear computational complexity with respect to the image size and thus directly scales to synthesizing high definition images. We show in the experiments that the proposed HiT achieves state-of-the-art FID scores of 30.83 and 2.95 on unconditional ImageNet \$128 \times 128\$ and FFHQ \$256 \times 256\$, respectively, with a reasonable throughput. We believe the proposed HiT is an important milestone for generators in GANs which are completely free of convolutions. Our code is made publicly available at <https://github.com/google-research/hit-gan>.

[Learning High-Precision Bounding Box for Rotated Object Detection via Kullback-Leibler Divergence](#)

- Xue Yang, Xiaojiang Yang, Jirui Yang, Qi Ming, Wentao Wang, Qi Tian, Junchi Yan
- abstract@[open-review](#): Existing rotated object detectors are mostly inherited from the horizontal detection paradigm, as the latter has evolved into a well-developed area. However, these detectors are difficult to perform prominently in high-precision detection due to the limitation of current regression loss design, especially for objects with large aspect ratios. Taking the perspective that horizontal detection is a special case for rotated object detection, in this paper, we are motivated to change the design of rotation regression loss from induction paradigm to deduction methodology, in terms of the relation between rotation and horizontal detection. We show that one essential challenge is how to modulate the coupled parameters in the rotation regression loss, as such the estimated parameters can influence to each other during the dynamic joint optimization, in an adaptive and synergetic way. Specifically, we first convert the rotated bounding box into a 2-D Gaussian distribution, and then calculate the Kullback-Leibler Divergence (KLD) between the Gaussian distributions as the regression loss. By analyzing the gradient of each parameter, we show that KLD (and its derivatives) can dynamically adjust the parameter gradients according to the characteristics of the object. For instance, it will adjust the importance (gradient weight) of the angle parameter according to the aspect ratio. This mechanism can be vital for high-precision detection as a slight angle error would cause a serious accuracy drop for large aspect ratios objects. More importantly, we have proved that KLD is scale invariant. We further show that the KLD loss can be degenerated into the popular Ln-norm loss for horizontal detection. Experimental results on seven datasets using different detectors show its consistent superiority, and codes are available at <https://github.com/yangxue0827/RotateDetection>.

[On Locality of Local Explanation Models](#)

- Sahra Ghalebikesabi, Lucile Ter-Minassian, Karla DiazOrdaz, Chris C Holmes
- abstract@[open-review](#): Shapley values provide model agnostic feature attributions for model outcome at a particular instance by simulating feature absence under a global population distribution. The use of a global population can lead to potentially misleading results when local model behaviour is of interest. Hence we consider the formulation of neighbourhood reference distributions that improve the local interpretability of Shapley values. By doing so, we find that the Nadaraya-Watson estimator, a well-studied kernel regressor, can be expressed as a self-normalised importance sampling estimator. Empirically, we observe that Neighbourhood Shapley values identify meaningful sparse feature relevance attributions that provide insight into local model behaviour, complimenting conventional Shapley analysis. They also increase on-manifold explainability and robustness to the construction of adversarial classifiers.

[FlexMatch: Boosting Semi-Supervised Learning with Curriculum Pseudo Labeling](#)

- Bowen Zhang, Yidong Wang, Wenxin Hou, HAO WU, Jindong Wang, Manabu Okumura, Takahiro Shinozaki
- abstract@[open-review](#): The recently proposed FixMatch achieved state-of-the-art results on most semi-supervised learning (SSL) benchmarks. However, like other modern SSL algorithms, FixMatch uses a pre-defined constant threshold for all classes to select unlabeled data that contribute to the training, thus failing to consider different learning status and learning difficulties of different classes. To address this issue, we propose Curriculum Pseudo Labeling (CPL), a curriculum learning approach to leverage unlabeled data according to the model's learning status. The core of CPL is to flexibly adjust thresholds for different classes at each time step to let pass informative unlabeled data and their pseudo labels. CPL does not introduce additional parameters or computations (forward or backward propagation). We apply CPL to FixMatch and call our improved algorithm FlexMatch. FlexMatch achieves state-of-the-art performance on a variety of SSL benchmarks, with especially strong performances when the labeled data are extremely limited or when the task is challenging. For example, FlexMatch achieves 13.96% and 18.96% error rate reduction over FixMatch on CIFAR-100 and STL-10 datasets respectively, when there are only 4 labels per class. CPL also significantly boosts the convergence speed, e.g., FlexMatch can use only 1/5 training time of FixMatch to achieve even better performance. Furthermore, we show that CPL can be easily adapted to other SSL algorithms and remarkably improve their performances. We open-source our code at <https://github.com/TorchSSL/TorchSSL>.

[Relative Flatness and Generalization](#)

- Henning Petzka, Michael Kamp, Linara Adilova, Cristian Sminchisescu, Mario Boley
- abstract@[open-review](#): Flatness of the loss curve is conjectured to be connected to the generalization ability of machine learning models, in particular neural networks. While it has been empirically observed that flatness measures consistently correlate strongly with generalization, it is still an open theoretical problem why and under which circumstances flatness is connected to generalization, in particular in light of reparameterizations that change certain flatness measures but leave generalization unchanged. We investigate the connection between flatness and generalization by relating it to the interpolation from representative data, deriving notions of representativeness, and feature robustness. The notions allow us to rigorously connect flatness and generalization and to identify conditions under which the connection holds. Moreover, they give rise to a novel, but natural relative flatness measure that correlates strongly with generalization, simplifies to ridge regression for ordinary least squares, and solves the reparameterization issue.

[The Image Local Autoregressive Transformer](#)

- Chenjie Cao, Yuxin Hong, Xiang Li, Chengrong Wang, Chengming Xu, Yanwei Fu, Xiangyang Xue
- abstract@[open-review](#): Recently, AutoRegressive (AR) models for the whole image generation empowered by transformers have achieved comparable or even better performance compared to Generative Adversarial Networks (GANs). Unfortunately, directly applying such AR models to edit/change local image regions, may suffer from the problems of missing global information, slow inference speed, and information leakage of local guidance. To address these limitations, we propose a novel model -- image Local Autoregressive Transformer (iLAT), to better facilitate the locally guided image synthesis. Our iLAT learns the novel local discrete representations, by the newly proposed local autoregressive (LA) transformer of the attention mask and convolution mechanism. Thus iLAT can efficiently synthesize the local image regions by key guidance information. Our iLAT is evaluated on various locally guided image syntheses, such as pose-guided person image synthesis and face editing. Both quantitative and qualitative results show the efficacy of our model.

[Towards Multi-Grained Explainability for Graph Neural Networks](#)

- Xiang Wang, Yingxin Wu, An Zhang, Xiangnan He, Tat-Seng Chua
- abstract@[open-review](#): When a graph neural network (GNN) made a prediction, one raises question about explainability: "Which fraction of the input graph is most influential to the model's decision?" Producing an answer requires understanding the model's inner workings in general and emphasizing the

insights on the decision for the instance at hand. Nonetheless, most of current approaches focus only on one aspect: (1) local explainability, which explains each instance independently, thus hardly exhibits the class-wise patterns; and (2) global explainability, which systematizes the globally important patterns, but might be trivial in the local context. This dichotomy limits the flexibility and effectiveness of explainers greatly. A performant paradigm towards multi-grained explainability is until-now lacking and thus a focus of our work. In this work, we exploit the pre-training and fine-tuning idea to develop our explainer and generate multi-grained explanations. Specifically, the pre-training phase accounts for the contrastivity among different classes, so as to highlight the class-wise characteristics from a global view; afterwards, the fine-tuning phase adapts the explanations in the local context. Experiments on both synthetic and real-world datasets show the superiority of our explainer, in terms of AUC on explaining graph classification over the leading baselines. Our codes and datasets are available at <https://github.com/Wuyxin/ReFine>.

[Behavior From the Void: Unsupervised Active Pre-Training](#)

- Hao Liu, Pieter Abbeel
- abstract@[open-review](#): We introduce a new unsupervised pre-training method for reinforcement learning called APT, which stands for Active Pre-Training. APT learns behaviors and representations by actively searching for novel states in reward-free environments. The key novel idea is to explore the environment by maximizing a non-parametric entropy computed in an abstract representation space, which avoids challenging density modeling and consequently allows our approach to scale much better in environments that have high-dimensional observations (e.g., image observations). We empirically evaluate APT by exposing task-specific reward after a long unsupervised pre-training phase. In Atari games, APT achieves human-level performance on 12 games and obtains highly competitive performance compared to canonical fully supervised RL algorithms. On DMControl suite, APT beats all baselines in terms of asymptotic performance and data efficiency and dramatically improves performance on tasks that are extremely difficult to train from scratch.

[Autonomous Reinforcement Learning via Subgoal Curricula](#)

- Archit Sharma, Abhishek Gupta, Sergey Levine, Karol Hausman, Chelsea Finn
- abstract@[open-review](#): Reinforcement learning (RL) promises to enable autonomous acquisition of complex behaviors for diverse agents. However, the success of current reinforcement learning algorithms is predicated on an often under-emphasised requirement -- each trial needs to start from a fixed initial state distribution. Unfortunately, resetting the environment to its initial state after each trial requires substantial amount of human supervision and extensive instrumentation of the environment which defeats the goal of autonomous acquisition of complex behaviors. In this work, we propose Value-accelerated Persistent Reinforcement Learning (VaPRL), which generates a curriculum of initial states such that the agent can bootstrap on the success of easier tasks to efficiently learn harder tasks. The agent also learns to reach the initial states proposed by the curriculum, minimizing the reliance on human interventions into the learning. We observe that VaPRL reduces the interventions required by three orders of magnitude compared to episodic RL while outperforming prior state-of-the art methods for reset-free RL both in terms of sample efficiency and asymptotic performance on a variety of simulated robotics problems.

[Statistically and Computationally Efficient Linear Meta-representation Learning](#)

- Kiran K. Thekumparampil, Prateek Jain, Praneeth Netrapalli, Sewoong Oh
- abstract@[open-review](#): In typical few-shot learning, each task is not equipped with enough data to be learned in isolation. To cope with such data scarcity, meta-representation learning methods train across many related tasks to find a shared (lower-dimensional) representation of the data where all tasks can be solved accurately. It is hypothesized that any new arriving tasks can be rapidly trained on this low-dimensional representation using only a few samples. Despite the practical successes of this approach, its statistical and computational properties are less understood. Moreover, the prescribed algorithms in these studies have little resemblance to those used in practice or they are computationally intractable. To understand and explain the success of popular meta-representation learning approaches such as ANIL, MetaOptNet, R2D2, and OML, we study a alternating gradient-descent minimization (AltMinGD) method (and its variant alternating minimization (AltMin)) which underlies the aforementioned methods. For a simple but canonical setting of shared linear representations, we show that AltMinGD achieves nearly-optimal estimation error, requiring only $\Omega(\mathrm{polylog}(d))$ samples per task. This agrees with the observed efficacy of this algorithm in the practical few-shot learning scenarios.

[Decentralized Learning in Online Queuing Systems](#)

- Flore Sentenac, Etienne Boursier, Vianney Perchet
- abstract@[open-review](#): Motivated by packet routing in computer networks, online queuing systems are composed of queues receiving packets at different rates. Repeatedly, they send packets to servers, each of them treating only at most one packet at a time. In the centralized case, the number of accumulated packets remains bounded (i.e., the system is stable) as long as the ratio between service rates and arrival rates is larger than \$1\$. In the decentralized case, individual no-regret strategies ensures stability when this ratio is larger than \$2\$. Yet, myopically minimizing regret disregards the long term effects due to the carryover of packets to further rounds. On the other hand, minimizing long term costs leads to stable Nash equilibria as soon as the ratio exceeds $\frac{e-1}{e}$. Stability with decentralized learning strategies with a ratio below \$2\$ was a major remaining question. We first argue that for ratios up to \$2\$, cooperation is required for stability of learning strategies, as selfish minimization of policy regret, a patient notion of regret, might indeed still be unstable in this case. We therefore consider cooperative queues and propose the first learning decentralized algorithm guaranteeing stability of the system as long as the ratio of rates is larger than \$1\$, thus reaching performances comparable to centralized strategies.

[Explainable Semantic Space by Grounding Language to Vision with Cross-Modal Contrastive Learning](#)

- Yizhen Zhang, Minkyu Choi, Kuan Han, Zhongming Liu
- abstract@[open-review](#): In natural language processing, most models try to learn semantic representations merely from texts. The learned representations encode the “distributional semantics” but fail to connect to any knowledge about the physical world. In contrast, humans learn language by grounding concepts in perception and action and the brain encodes “grounded semantics” for cognition. Inspired by this notion and recent work in vision-language learning, we design a two-stream model for grounding language learning in vision. The model includes a VGG-based visual stream and a Bert-based language stream. The two streams merge into a joint representational space. Through cross-modal contrastive learning, the model first learns to align visual and language representations with the MS COCO dataset. The model further learns to retrieve visual objects with language queries through a cross-modal attention module and to infer the visual relations between the retrieved objects through a bilinear operator with the Visual Genome dataset. After training, the model’s language stream is a stand-alone language model capable of embedding concepts in a visually grounded semantic space. This semantic space manifests principal dimensions explainable with human intuition and neurobiological knowledge. Word embeddings in this semantic space are predictive of human-defined norms of semantic features and are segregated into perceptually distinctive clusters. Furthermore, the visually grounded language model also enables compositional language understanding based on visual knowledge and multimodal image search with queries based on images, texts, or their combinations.

[BulletTrain: Accelerating Robust Neural Network Training via Boundary Example Mining](#)

- Weizhe Hua, Yichi Zhang, Chuan Guo, Zhiru Zhang, G. Edward Suh
- abstract@[open-review](#): Neural network robustness has become a central topic in machine learning in recent years. Most training algorithms that improve the model’s robustness to adversarial and common corruptions also introduce a large computational overhead, requiring as many as ten times the number of forward and backward passes in order to converge. To combat this inefficiency, we propose BulletTrain, a boundary example mining technique to

drastically reduce the computational cost of robust training. Our key observation is that only a small fraction of examples are beneficial for improving robustness. BulletTrain dynamically predicts these important examples and optimizes robust training algorithms to focus on the important examples. We apply our technique to several existing robust training algorithms and achieve a 2.2x speed-up for TRADES and MART on CIFAR-10 and a 1.7x speed-up for AugMix on CIFAR-10-C and CIFAR-100-C without any reduction in clean and robust accuracy.

[Neural Distance Embeddings for Biological Sequences](#)

- Gabriele Corso, Zhitao Ying, Michal Pády, Petar Veličković, Jure Leskovec, Pietro Liò
- abstract@[open-review](#): The development of data-dependent heuristics and representations for biological sequences that reflect their evolutionary distance is critical for large-scale biological research. However, popular machine learning approaches, based on continuous Euclidean spaces, have struggled with the discrete combinatorial formulation of the edit distance that models evolution and the hierarchical relationship that characterises real-world datasets. We present Neural Distance Embeddings (NeuroSEED), a general framework to embed sequences in geometric vector spaces, and illustrate the effectiveness of the hyperbolic space that captures the hierarchical structure and provides an average 38% reduction in embedding RMSE against the best competing geometry. The capacity of the framework and the significance of these improvements are then demonstrated devising supervised and unsupervised NeuroSEED approaches to multiple core tasks in bioinformatics. Benchmarked with common baselines, the proposed approaches display significant accuracy and/or runtime improvements on real-world datasets. As an example for hierarchical clustering, the proposed pretrained and from-scratch methods match the quality of competing baselines with 30x and 15x runtime reduction, respectively.

[Fitting summary statistics of neural data with a differentiable spiking network simulator](#)

- Guillaume Bellec, Shuqi Wang, Alireza Modirshanechi, Johann Brea, Wulfram Gerstner
- abstract@[open-review](#): Fitting network models to neural activity is an important tool in neuroscience. A popular approach is to model a brain area with a probabilistic recurrent spiking network whose parameters maximize the likelihood of the recorded activity. Although this is widely used, we show that the resulting model does not produce realistic neural activity. To correct for this, we suggest to augment the log-likelihood with terms that measure the dissimilarity between simulated and recorded activity. This dissimilarity is defined via summary statistics commonly used in neuroscience and the optimization is efficient because it relies on back-propagation through the stochastically simulated spike trains. We analyze this method theoretically and show empirically that it generates more realistic activity statistics. We find that it improves upon other fitting algorithms for spiking network models like GLMs (Generalized Linear Models) which do not usually rely on back-propagation. This new fitting algorithm also enables the consideration of hidden neurons which is otherwise notoriously hard, and we show that it can be crucial when trying to infer the network connectivity from spike recordings.

[PerSim: Data-Efficient Offline Reinforcement Learning with Heterogeneous Agents via Personalized Simulators](#)

- Anish Agarwal, Abdullah Alomar, Varkey Alumootil, Devavrat Shah, Dennis Shen, Zhi Xu, Cindy Yang
- abstract@[open-review](#): We consider offline reinforcement learning (RL) with heterogeneous agents under severe data scarcity, i.e., we only observe a single historical trajectory for every agent under an unknown, potentially sub-optimal policy. We find that the performance of state-of-the-art offline and model-based RL methods degrade significantly given such limited data availability, even for commonly perceived "solved" benchmark settings such as "MountainCar" and "CartPole". To address this challenge, we propose PerSim, a model-based offline RL approach which first learns a personalized simulator for each agent by collectively using the historical trajectories across all agents, prior to learning a policy. We do so by positing that the transition dynamics across agents can be represented as a latent function of latent factors associated with agents, states, and actions; subsequently, we theoretically establish that this function is well-approximated by a "low-rank" decomposition of separable agent, state, and action latent functions. This representation suggests a simple, regularized neural network architecture to effectively learn the transition dynamics per agent, even with scarce, offline data. We perform extensive experiments across several benchmark environments and RL methods. The consistent improvement of our approach, measured in terms of both state dynamics prediction and eventual reward, confirms the efficacy of our framework in leveraging limited historical data to simultaneously learn personalized policies across agents.

[Online Sign Identification: Minimization of the Number of Errors in Thresholding Bandits](#)

- Reda Ouhamma, Odalric-Ambrym Maillard, Vianney Perchet
- abstract@[open-review](#): In the fixed budget thresholding bandit problem, an algorithm sequentially allocates a budgeted number of samples to different distributions. It then predicts whether the mean of each distribution is larger or lower than a given threshold. We introduce a large family of algorithms (containing most existing relevant ones), inspired by the Frank-Wolfe algorithm, and provide a thorough yet generic analysis of their performance. This allowed us to construct new explicit algorithms, for a broad class of problems, whose losses are within a small constant factor of the non-adaptive oracle ones. Quite interestingly, we observed that adaptive methods empirically greatly out-perform non-adaptive oracles, an uncommon behavior in standard online learning settings, such as regret minimization. We explain this surprising phenomenon on an insightful toy problem.

[All Tokens Matter: Token Labeling for Training Better Vision Transformers](#)

- Zi-Hang Jiang, Qibin Hou, Li Yuan, Daquan Zhou, Yujun Shi, Xiaojie Jin, Anran Wang, Jiashi Feng
- abstract@[open-review](#): In this paper, we present token labeling---a new training objective for training high-performance vision transformers (ViTs). Different from the standard training objective of ViTs that computes the classification loss on an additional trainable class token, our proposed one takes advantage of all the image patch tokens to compute the training loss in a dense manner. Specifically, token labeling reformulates the image classification problem into multiple token-level recognition problems and assigns each patch token with an individual location-specific supervision generated by a machine annotator. Experiments show that token labeling can clearly and consistently improve the performance of various ViT models across a wide spectrum. For a vision transformer with 26M learnable parameters serving as an example, with token labeling, the model can achieve 84.4% Top-1 accuracy on ImageNet. The result can be further increased to 86.4% by slightly scaling the model size up to 150M, delivering the minimal-sized model among previous models (250M+) reaching 86%. We also show that token labeling can clearly improve the generalization capability of the pretrained models on downstream tasks with dense prediction, such as semantic segmentation. Our code and model are publicly available at <https://github.com/zihangJiang/TokenLabeling>.

[Partition and Code: learning how to compress graphs](#)

- Giorgos Bouritsas, Andreas Loukas, Nikolaos Karalias, Michael Bronstein
- abstract@[open-review](#): Can we use machine learning to compress graph data? The absence of ordering in graphs poses a significant challenge to conventional compression algorithms, limiting their attainable gains as well as their ability to discover relevant patterns. On the other hand, most graph compression approaches rely on domain-dependent handcrafted representations and cannot adapt to different underlying graph distributions. This work aims to establish the necessary principles a lossless graph compression method should follow to approach the entropy storage lower bound. Instead of making rigid assumptions about the graph distribution, we formulate the compressor as a probabilistic model that can be learned from data and generalise to unseen instances. Our "Partition and Code" framework entails three steps: first, a partitioning algorithm decomposes the graph into subgraphs, then these are mapped to the elements of a small dictionary on which we learn a probability distribution, and finally, an entropy encoder translates the representation into bits. All the components (partitioning, dictionary and distribution) are parametric and can be trained with gradient descent. We theoretically compare the compression quality of several graph encodings and prove, under mild conditions, that PnC achieves compression gains that

grow either linearly or quadratically with the number of vertices. Empirically, PnC yields significant compression improvements on diverse real-world networks.

[Knowledge-inspired 3D Scene Graph Prediction in Point Cloud](#)

- Shoulong Zhang, shuai li, Aimin Hao, Hong Qin
- abstract@[open-review](#): Prior knowledge integration helps identify semantic entities and their relationships in a graphical representation, however, its meaningful abstraction and intervention remain elusive. This paper advocates a knowledge-inspired 3D scene graph prediction method solely based on point clouds. At the mathematical modeling level, we formulate the task as two sub-problems: knowledge learning and scene graph prediction with learned prior knowledge. Unlike conventional methods that learn knowledge embedding and regular patterns from encoded visual information, we propose to suppress the misunderstandings caused by appearance similarities and other perceptual confusion. At the network design level, we devise a graph auto-encoder to automatically extract class-dependent representations and topological patterns from the one-hot class labels and their intrinsic graphical structures, so that the prior knowledge can avoid perceptual errors and noises. We further devise a scene graph prediction model to predict credible relationship triplets by incorporating the related prototype knowledge with perceptual information. Comprehensive experiments confirm that, our method can successfully learn representative knowledge embedding, and the obtained prior knowledge can effectively enhance the accuracy of relationship predictions. Our thorough evaluations indicate the new method can achieve the state-of-the-art performance compared with other scene graph prediction methods.

[Online Variational Filtering and Parameter Learning](#)

- Andrew Campbell, Yuyang Shi, Thomas Rainforth, Arnaud Doucet
- abstract@[open-review](#): We present a variational method for online state estimation and parameter learning in state-space models (SSMs), a ubiquitous class of latent variable models for sequential data. As per standard batch variational techniques, we use stochastic gradients to simultaneously optimize a lower bound on the log evidence with respect to both model parameters and a variational approximation of the states' posterior distribution. However, unlike existing approaches, our method is able to operate in an entirely online manner, such that historic observations do not require revisit after being incorporated and the cost of updates at each time step remains constant, despite the growing dimensionality of the joint posterior distribution of the states. This is achieved by utilizing backward decompositions of this joint posterior distribution and of its variational approximation, combined with Bellman-type recursions for the evidence lower bound and its gradients. We demonstrate the performance of this methodology across several examples, including high-dimensional SSMs and sequential Variational Auto-Encoders.

[Heavy Ball Neural Ordinary Differential Equations](#)

- Hedi Xia, Vai Suliafu, Hangjie Ji, Tan Nguyen, Andrea Bertozzi, Stanley Osher, Bao Wang
- abstract@[open-review](#): We propose heavy ball neural ordinary differential equations (HBNODEs), leveraging the continuous limit of the classical momentum accelerated gradient descent, to improve neural ODEs (NODEs) training and inference. HBNODEs have two properties that imply practical advantages over NODEs: (i) The adjoint state of an HBNODE also satisfies an HBNODE, accelerating both forward and backward ODE solvers, thus significantly reducing the number of function evaluations (NFEs) and improving the utility of the trained models. (ii) The spectrum of HBNODEs is well structured, enabling effective learning of long-term dependencies from complex sequential data. We verify the advantages of HBNODEs over NODEs on benchmark tasks, including image classification, learning complex dynamics, and sequential modeling. Our method requires remarkably fewer forward and backward NFEs, is more accurate, and learns long-term dependencies more effectively than the other ODE-based neural network models. Code is available at \url{<https://github.com/hedixia/HeavyBallNODE>}.

[Structure learning in polynomial time: Greedy algorithms, Bregman information, and exponential families](#)

- Goutham Rajendran, Bohdan Kivva, Ming Gao, Bryon Aragam
- abstract@[open-review](#): Greedy algorithms have long been a workhorse for learning graphical models, and more broadly for learning statistical models with sparse structure. In the context of learning directed acyclic graphs, greedy algorithms are popular despite their worst-case exponential runtime. In practice, however, they are very efficient. We provide new insight into this phenomenon by studying a general greedy score-based algorithm for learning DAGs. Unlike edge-greedy algorithms such as the popular GES and hill-climbing algorithms, our approach is vertex-greedy and requires at most a polynomial number of score evaluations. We then show how recent polynomial-time algorithms for learning DAG models are a special case of this algorithm, thereby illustrating how these order-based algorithms can be rigorously interpreted as score-based algorithms. This observation suggests new score functions and optimality conditions based on the duality between Bregman divergences and exponential families, which we explore in detail. Explicit sample and computational complexity bounds are derived. Finally, we provide extensive experiments suggesting that this algorithm indeed optimizes the score in a variety of settings.

[On the Sample Complexity of Learning under Geometric Stability](#)

- Alberto Bietti, Luca Venturi, Joan Bruna
- abstract@[open-review](#): Many supervised learning problems involve high-dimensional data such as images, text, or graphs. In order to make efficient use of data, it is often useful to leverage certain geometric priors in the problem at hand, such as invariance to translations, permutation subgroups, or stability to small deformations. We study the sample complexity of learning problems where the target function presents such invariance and stability properties, by considering spherical harmonic decompositions of such functions on the sphere. We provide non-parametric rates of convergence for kernel methods, and show improvements in sample complexity by a factor equal to the size of the group when using an invariant kernel over the group, compared to the corresponding non-invariant kernel. These improvements are valid when the sample size is large enough, with an asymptotic behavior that depends on spectral properties of the group. Finally, these gains are extended beyond invariance groups to also cover geometric stability to small deformations, modeled here as subsets (not necessarily subgroups) of permutations.

[SIMILAR: Submodular Information Measures Based Active Learning In Realistic Scenarios](#)

- Suraj Kothawade, Nathan Beck, Krishnateja Killamsetty, Rishabh Iyer
- abstract@[open-review](#): Active learning has proven to be useful for minimizing labeling costs by selecting the most informative samples. However, existing active learning methods do not work well in realistic scenarios such as imbalance or rare classes, out-of-distribution data in the unlabeled set, and redundancy. In this work, we propose SIMILAR (Submodular Information Measures based actIve LeARning), a unified active learning framework using recently proposed submodular information measures (SIM) as acquisition functions. We argue that SIMILAR not only works in standard active learning but also easily extends to the realistic settings considered above and acts as a one-stop solution for active learning that is scalable to large real-world datasets. Empirically, we show that SIMILAR significantly outperforms existing active learning algorithms by as much as ~5%–18% in the case of rare classes and ~5%–10% in the case of out-of-distribution data on several image classification tasks like CIFAR-10, MNIST, and ImageNet.

[Monte Carlo Tree Search With Iteratively Refining State Abstractions](#)

- Samuel Sokota, Caleb Y Ho, Zaheen Ahmad, J. Zico Kolter

- abstract@[open-review](#): Decision-time planning is the process of constructing a transient, local policy with the intent of using it to make the immediate decision. Monte Carlo tree search (MCTS), which has been leveraged to great success in Go, chess, shogi, Hex, Atari, and other settings, is perhaps the most celebrated decision-time planning algorithm. Unfortunately, in its original form, MCTS can degenerate to one-step search in domains with stochasticity. Progressive widening is one way to ameliorate this issue, but we argue that it possesses undesirable properties for some settings. In this work, we present a method, called abstraction refining, for extending MCTS to stochastic environments which, unlike progressive widening, leverages the geometry of the state space. We argue that leveraging the geometry of the space can offer advantages. To support this claim, we present a series of experimental examples in which abstraction refining outperforms progressive widening, given equal simulation budgets.

[Flattening Sharpness for Dynamic Gradient Projection Memory Benefits Continual Learning](#)

- Danruo DENG, Guangyong Chen, Jianye Hao, Qiong Wang, Pheng-Ann Heng
- abstract@[open-review](#): The backpropagation networks are notably susceptible to catastrophic forgetting, where networks tend to forget previously learned skills upon learning new ones. To address such the 'sensitivity-stability' dilemma, most previous efforts have been contributed to minimizing the empirical risk with different parameter regularization terms and episodic memory, but rarely exploring the usages of the weight loss landscape. In this paper, we investigate the relationship between the weight loss landscape and sensitivity-stability in the continual learning scenario, based on which, we propose a novel method, Flattening Sharpness for Dynamic Gradient Projection Memory (FS-DGPM). In particular, we introduce a soft weight to represent the importance of each basis representing past tasks in GPM, which can be adaptively learned during the learning process, so that less important bases can be dynamically released to improve the sensitivity of new skill learning. We further introduce Flattening Sharpness (FS) to reduce the generalization gap by explicitly regulating the flatness of the weight loss landscape of all seen tasks. As demonstrated empirically, our proposed method consistently outperforms baselines with the superior ability to learn new skills while alleviating forgetting effectively.

[Taxonomizing local versus global structure in neural network loss landscapes](#)

- Yaoqing Yang, Liam Hodgkinson, Ryan Theisen, Joe Zou, Joseph E. Gonzalez, Kannan Ramchandran, Michael W. Mahoney
- abstract@[open-review](#): Viewing neural network models in terms of their loss landscapes has a long history in the statistical mechanics approach to learning, and in recent years it has received attention within machine learning proper. Among other things, local metrics (such as the smoothness of the loss landscape) have been shown to correlate with global properties of the model (such as good generalization performance). Here, we perform a detailed empirical analysis of the loss landscape structure of thousands of neural network models, systematically varying learning tasks, model architectures, and/or quantity/quality of data. By considering a range of metrics that attempt to capture different aspects of the loss landscape, we demonstrate that the best test accuracy is obtained when: the loss landscape is globally well-connected; ensembles of trained models are more similar to each other; and models converge to locally smooth regions. We also show that globally poorly-connected landscapes can arise when models are small or when they are trained to lower quality data; and that, if the loss landscape is globally poorly-connected, then training to zero loss can actually lead to worse test accuracy. Our detailed empirical results shed light on phases of learning (and consequent double descent behavior), fundamental versus incidental determinants of good generalization, the role of load-like and temperature-like parameters in the learning process, different influences on the loss landscape from model and data, and the relationships between local and global metrics, all topics of recent interest.

[Learning Models for Actionable Recourse](#)

- Alexis Ross, Himabindu Lakkaraju, Osbert Bastani
- abstract@[open-review](#): As machine learning models are increasingly deployed in high-stakes domains such as legal and financial decision-making, there has been growing interest in post-hoc methods for generating counterfactual explanations. Such explanations provide individuals adversely impacted by predicted outcomes (e.g., an applicant denied a loan) with recourse---i.e., a description of how they can change their features to obtain a positive outcome. We propose a novel algorithm that leverages adversarial training and PAC confidence sets to learn models that theoretically guarantee recourse to affected individuals with high probability without sacrificing accuracy. We demonstrate the efficacy of our approach via extensive experiments on real data.

[Efficient and Accurate Gradients for Neural SDEs](#)

- Patrick Kidger, James Foster, Xuechen (Chen) Li, Terry Lyons
- abstract@[open-review](#): Neural SDEs combine many of the best qualities of both RNNs and SDEs, and as such are a natural choice for modelling many types of temporal dynamics. They offer memory efficiency, high-capacity function approximation, and strong priors on model space. Neural SDEs may be trained as VAEs or as GANs; in either case it is necessary to backpropagate through the SDE solve. In particular this may be done by constructing a backwards-in-time SDE whose solution is the desired parameter gradients. However, this has previously suffered from severe speed and accuracy issues, due to high computational complexity, numerical errors in the SDE solve, and the cost of reconstructing Brownian motion. Here, we make several technical innovations to overcome these issues. First, we introduce the \textit{reversible Heun method}: a new SDE solver that is algebraically reversible --- which reduces numerical gradient errors to almost zero, improving several test metrics by substantial margins over state-of-the-art. Moreover it requires half as many function evaluations as comparable solvers, giving up to a \$1.98\times\$ speedup. Next, we introduce the \textit{Brownian interval}. This is a new and computationally efficient way of exactly sampling \textit{and reconstructing} Brownian motion; this is in contrast to previous reconstruction techniques that are both approximate and relatively slow. This gives up to a \$10.6\times\$ speed improvement over previous techniques. After that, when specifically training Neural SDEs as GANs (Kidger et al. 2021), we demonstrate how SDE-GANs may be trained through careful weight clipping and choice of activation function. This reduces computational cost (giving up to a \$1.87\times\$ speedup), and removes the truncation errors of the double adjoint required for gradient penalty, substantially improving several test metrics. Altogether these techniques offer substantial improvements over the state-of-the-art, with respect to both training speed and with respect to classification, prediction, and MMD test metrics. We have contributed implementations of all of our techniques to the \texttt{torchsde} library to help facilitate their adoption.

[EIGNN: Efficient Infinite-Depth Graph Neural Networks](#)

- Juncheng Liu, Kenji Kawaguchi, Bryan Hooi, Yiwei Wang, Xiaokui Xiao
- abstract@[open-review](#): Graph neural networks (GNNs) are widely used for modelling graph-structured data in numerous applications. However, with their inherently finite aggregation layers, existing GNN models may not be able to effectively capture long-range dependencies in the underlying graphs. Motivated by this limitation, we propose a GNN model with infinite depth, which we call Efficient Infinite-Depth Graph Neural Networks (EIGNN), to efficiently capture very long-range dependencies. We theoretically derive a closed-form solution of EIGNN which makes training an infinite-depth GNN model tractable. We then further show that we can achieve more efficient computation for training EIGNN by using eigendecomposition. The empirical results of comprehensive experiments on synthetic and real-world datasets show that EIGNN has a better ability to capture long-range dependencies than recent baselines, and consistently achieves state-of-the-art performance. Furthermore, we show that our model is also more robust against both noise and adversarial perturbations on node features.

[Fractal Structure and Generalization Properties of Stochastic Optimization Algorithms](#)

- Alexander Camuto, George Deligiannidis, Murat A. Erdogdu, Mert Gurbuzbalaban, Umut Simsekli, Lingjiong Zhu
- abstract@[open-review](#): Understanding generalization in deep learning has been one of the major challenges in statistical learning theory over the last decade. While recent work has illustrated that the dataset and the training algorithm must be taken into account in order to obtain meaningful generalization bounds, it is still theoretically not clear which properties of the data and the algorithm determine the generalization performance. In this

study, we approach this problem from a dynamical systems theory perspective and represent stochastic optimization algorithms as random iterated function systems (IFS). Well studied in the dynamical systems literature, under mild assumptions, such IFSs can be shown to be ergodic with an invariant measure that is often supported on sets with a fractal structure. As our main contribution, we prove that the generalization error of a stochastic optimization algorithm can be bounded based on the ‘complexity’ of the fractal structure that underlies its invariant measure. Then, by leveraging results from dynamical systems theory, we show that the generalization error can be explicitly linked to the choice of the algorithm (e.g., stochastic gradient descent -- SGD), algorithm hyperparameters (e.g., step-size, batch-size), and the geometry of the problem (e.g., Hessian of the loss). We further specialize our results to specific problems (e.g., linear/logistic regression, one hidden-layered neural networks) and algorithms (e.g., SGD and preconditioned variants), and obtain analytical estimates for our bound. For modern neural networks, we develop an efficient algorithm to compute the developed bound and support our theory with various experiments on neural networks.

[An Infinite-Feature Extension for Bayesian ReLU Nets That Fixes Their Asymptotic Overconfidence](#)

- Agustinus Kristiadi, Matthias Hein, Philipp Hennig
- abstract@[open-review](#): A Bayesian treatment can mitigate overconfidence in ReLU nets around the training data. But far away from them, ReLU Bayesian neural networks (BNNs) can still underestimate uncertainty and thus be asymptotically overconfident. This issue arises since the output variance of a BNN with finitely many features is quadratic in the distance from the data region. Meanwhile, Bayesian linear models with ReLU features converge, in the infinite-width limit, to a particular Gaussian process (GP) with a variance that grows cubically so that no asymptotic overconfidence can occur. While this may seem of mostly theoretical interest, in this work, we show that it can be used in practice to the benefit of BNNs. We extend finite ReLU BNNs with infinite ReLU features via the GP and show that the resulting model is asymptotically maximally uncertain far away from the data while the BNNs' predictive power is unaffected near the data. Although the resulting model approximates a full GP posterior, thanks to its structure, it can be applied post-hoc to any pre-trained ReLU BNN at a low cost.

[Bandit Phase Retrieval](#)

- Tor Lattimore, Botao Hao
- abstract@[open-review](#): We study a bandit version of phase retrieval where the learner chooses actions $(A_t)_{t=1}^n$ in the d -dimensional unit ball and the expected reward is $\langle A_t, \theta^* \rangle^2$ with $\theta^* \in \mathbb{R}^d$ an unknown parameter vector. We prove an upper bound on the minimax cumulative regret in this problem of $\tilde{\Theta}(d \sqrt{n})$, which matches known lower bounds up to logarithmic factors and improves on the best known upper bound by a factor of \sqrt{d} . We also show that the minimax simple regret is $\tilde{\Theta}(d / \sqrt{n})$ and that this is only achievable by an adaptive algorithm. Our analysis shows that an apparently convincing heuristic for guessing lower bounds can be misleading and that uniform bounds on the information ratio for information-directed sampling (Russo and Van Roy, 2014) are not sufficient for optimal regret.

[Lower Bounds on Metropolized Sampling Methods for Well-Conditioned Distributions](#)

- Yin Tat Lee, Ruoqi Shen, Kevin Tian
- abstract@[open-review](#): We give lower bounds on the performance of two of the most popular sampling methods in practice, the Metropolis-adjusted Langevin algorithm (MALA) and multi-step Hamiltonian Monte Carlo (HMC) with a leapfrog integrator, when applied to well-conditioned distributions. Our main result is a nearly-tight lower bound of $\widetilde{\Omega}(\kappa d)$ on the mixing time of MALA from an exponentially warm start, matching a line of algorithmic results [cite{DwivediCW018, ChenDWY19, LeeST20a}](#) up to logarithmic factors and answering an open question of [cite{ChewiLACGR20}](#). We also show that a polynomial dependence on dimension is necessary for the relaxation time of HMC under any number of leapfrog steps, and bound the gains achievable by changing the step count. Our HMC analysis draws upon a novel connection between leapfrog integration and Chebyshev polynomials, which may be of independent interest.

[Taming Communication and Sample Complexities in Decentralized Policy Evaluation for Cooperative Multi-Agent Reinforcement Learning](#)

- Xin Zhang, Zhuqing Liu, Jia Liu, Zhengyuan Zhu, Songtao Lu
- abstract@[open-review](#): Cooperative multi-agent reinforcement learning (MARL) has received increasing attention in recent years and has found many scientific and engineering applications. However, a key challenge arising from many cooperative MARL algorithm designs (e.g., the actor-critic framework) is the policy evaluation problem, which can only be conducted in a decentralized fashion. In this paper, we focus on decentralized MARL policy evaluation with nonlinear function approximation, which is often seen in deep MARL. We first show that the empirical decentralized MARL policy evaluation problem can be reformulated as a decentralized nonconvex-strongly-concave minimax saddle point problem. We then develop a decentralized gradient-based descent ascent algorithm called GT-GDA that enjoys a convergence rate of $\mathcal{O}(1/T)$. To further reduce the sample complexity, we propose two decentralized stochastic optimization algorithms called GT-SRVR and GT-SRVRI, which enhance GT-GDA by variance reduction techniques. We show that all algorithms all enjoy an $\mathcal{O}(1/T)$ convergence rate to a stationary point of the reformulated minimax problem. Moreover, the fast convergence rates of GT-SRVR and GT-SRVRI imply $\mathcal{O}(\epsilon^{-2})$ communication complexity and $\mathcal{O}(m\sqrt{n}\epsilon^{-2})$ sample complexity, where m is the number of agents and n is the length of trajectories. To our knowledge, this paper is the first work that achieves both $\mathcal{O}(\epsilon^{-2})$ sample complexity and $\mathcal{O}(\epsilon^{-2})$ communication complexity in decentralized policy evaluation for cooperative MARL. Our extensive experiments also corroborate the theoretical performance of our proposed decentralized policy evaluation algorithms.

[Federated Graph Classification over Non-IID Graphs](#)

- Han Xie, Jing Ma, Li Xiong, Carl Yang
- abstract@[open-review](#): Federated learning has emerged as an important paradigm for training machine learning models in different domains. For graph-level tasks such as graph classification, graphs can also be regarded as a special type of data samples, which can be collected and stored in separate local systems. Similar to other domains, multiple local systems, each holding a small set of graphs, may benefit from collaboratively training a powerful graph mining model, such as the popular graph neural networks (GNNs). To provide more motivation towards such endeavors, we analyze real-world graphs from different domains to confirm that they indeed share certain graph properties that are statistically significant compared with random graphs. However, we also find that different sets of graphs, even from the same domain or same dataset, are non-IID regarding both graph structures and node features. To handle this, we propose a graph clustered federated learning (GCFL) framework that dynamically finds clusters of local systems based on the gradients of GNNs, and theoretically justify that such clusters can reduce the structure and feature heterogeneity among graphs owned by the local systems. Moreover, we observe the gradients of GNNs to be rather fluctuating in GCFL which impedes high-quality clustering, and design a gradient sequence-based clustering mechanism based on dynamic time warping (GCFL+). Extensive experimental results and in-depth analysis demonstrate the effectiveness of our proposed frameworks.

[SubTab: Subsetting Features of Tabular Data for Self-Supervised Representation Learning](#)

- Talip Ucar, Ehsan Hajiramezanali, Lindsay Edwards
- abstract@[open-review](#): Self-supervised learning has been shown to be very effective in learning useful representations, and yet much of the success is achieved in data types such as images, audio, and text. The success is mainly enabled by taking advantage of spatial, temporal, or semantic structure in the

data through augmentation. However, such structure may not exist in tabular datasets commonly used in fields such as healthcare, making it difficult to design an effective augmentation method, and hindering a similar progress in tabular data setting. In this paper, we introduce a new framework, Subsetting features of Tabular data (SubTab), that turns the task of learning from tabular data into a multi-view representation learning problem by dividing the input features to multiple subsets. We argue that reconstructing the data from the subset of its features rather than its corrupted version in an autoencoder setting can better capture its underlying latent representation. In this framework, the joint representation can be expressed as the aggregate of latent variables of the subsets at test time, which we refer to as collaborative inference. Our experiments show that the SubTab achieves the state of the art (SOTA) performance of 98.31% on MNIST in tabular setting, on par with CNN-based SOTA models, and surpasses existing baselines on three other real-world datasets by a significant margin.

Convergence Rates of Stochastic Gradient Descent under Infinite Noise Variance

- Hongjian Wang, Mert Gurbuzbalaban, Lingjiong Zhu, Umut Simsekli, Murat A. Erdogdu
- abstract@[open-review](#): Recent studies have provided both empirical and theoretical evidence illustrating that heavy tails can emerge in stochastic gradient descent (SGD) in various scenarios. Such heavy tails potentially result in iterates with diverging variance, which hinders the use of conventional convergence analysis techniques that rely on the existence of the second-order moments. In this paper, we provide convergence guarantees for SGD under a state-dependent and heavy-tailed noise with a potentially infinite variance, for a class of strongly convex objectives. In the case where the $\$p\$$ -th moment of the noise exists for some $\$p \in [1,2]\$$, we first identify a condition on the Hessian, coined ' $\$p\$$ -positive (semi-)definiteness', that leads to an interesting interpolation between the positive semi-definite cone ($\$p=2\$\$$) and the cone of diagonally dominant matrices with non-negative diagonal entries ($\$p=1\$\$$). Under this condition, we provide a convergence rate for the distance to the global optimum in $\$L^p\$$. Furthermore, we provide a generalized central limit theorem, which shows that the properly scaled Polyak-Ruppert averaging converges weakly to a multivariate $\$alpha\$$ -stable random vector. Our results indicate that even under heavy-tailed noise with infinite variance, SGD can converge to the global optimum without necessitating any modification neither to the loss function nor to the algorithm itself, as typically required in robust statistics. We demonstrate the implications of our results over misspecified models, in the presence of heavy-tailed data.

Conflict-Averse Gradient Descent for Multi-task learning

- Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, Qiang Liu
- abstract@[open-review](#): The goal of multi-task learning is to enable more efficient learning than single task learning by sharing model structures for a diverse set of tasks. A standard multi-task learning objective is to minimize the average loss across all tasks. While straightforward, using this objective often results in much worse final performance for each task than learning them independently. A major challenge in optimizing a multi-task model is the conflicting gradients, where gradients of different task objectives are not well aligned so that following the average gradient direction can be detrimental to specific tasks' performance. Previous work has proposed several heuristics to manipulate the task gradients for mitigating this problem. But most of them lack convergence guarantee and/or could converge to any Pareto-stationary point. In this paper, we introduce Conflict-Averse Gradient descent (CAGrad) which minimizes the average loss function, while leveraging the worst local improvement of individual tasks to regularize the algorithm trajectory. CAGrad balances the objectives automatically and still provably converges to a minimum over the average loss. It includes the regular gradient descent (GD) and the multiple gradient descent algorithm (MGDA) in the multi-objective optimization (MOO) literature as special cases. On a series of challenging multi-task supervised learning and reinforcement learning tasks, CAGrad achieves improved performance over prior state-of-the-art multi-objective gradient manipulation methods.

Amortized Synthesis of Constrained Configurations Using a Differentiable Surrogate

- Xingyuan Sun, Tianju Xue, Szymon Rusinkiewicz, Ryan P. Adams
- abstract@[open-review](#): In design, fabrication, and control problems, we are often faced with the task of synthesis, in which we must generate an object or configuration that satisfies a set of constraints while maximizing one or more objective functions. The synthesis problem is typically characterized by a physical process in which many different realizations may achieve the goal. This many-to-one map presents challenges to the supervised learning of feed-forward synthesis, as the set of viable designs may have a complex structure. In addition, the non-differentiable nature of many physical simulations prevents efficient direct optimization. We address both of these problems with a two-stage neural network architecture that we may consider to be an autoencoder. We first learn the decoder: a differentiable surrogate that approximates the many-to-one physical realization process. We then learn the encoder, which maps from goal to design, while using the fixed decoder to evaluate the quality of the realization. We evaluate the approach on two case studies: extruder path planning in additive manufacturing and constrained soft robot inverse kinematics. We compare our approach to direct optimization of the design using the learned surrogate, and to supervised learning of the synthesis problem. We find that our approach produces higher quality solutions than supervised learning, while being competitive in quality with direct optimization, at a greatly reduced computational cost.

Efficient First-Order Contextual Bandits: Prediction, Allocation, and Triangular Discrimination

- Dylan J. Foster, Akshay Krishnamurthy
- abstract@[open-review](#): A recurring theme in statistical learning, online learning, and beyond is that faster convergence rates are possible for problems with low noise, often quantified by the performance of the best hypothesis; such results are known as first-order or small-loss guarantees. While first-order guarantees are relatively well understood in statistical and online learning, adapting to low noise in contextual bandits (and more broadly, decision making) presents major algorithmic challenges. In a COLT 2017 open problem, Agarwal, Krishnamurthy, Langford, Luo, and Schapire asked whether first-order guarantees are even possible for contextual bandits and---if so---whether they can be attained by efficient algorithms. We give a resolution to this question by providing an optimal and efficient reduction from contextual bandits to online regression with the logarithmic (or, cross-entropy) loss. Our algorithm is simple and practical, readily accommodates rich function classes, and requires no distributional assumptions beyond realizability. In a large-scale empirical evaluation, we find that our approach typically outperforms comparable non-first-order methods. On the technical side, we show that the logarithmic loss and an information-theoretic quantity called the triangular discrimination play a fundamental role in obtaining first-order guarantees, and we combine this observation with new refinements to the regression oracle reduction framework of Foster and Rakhlin (2020). The use of triangular discrimination yields novel results even for the classical statistical learning model, and we anticipate that it will find broader use.

Distributed Estimation with Multiple Samples per User: Sharp Rates and Phase Transition

- Jayadev Acharya, Clement Canonne, Yuhan Liu, Ziteng Sun, Himanshu Tyagi
- abstract@[open-review](#): We obtain tight minimax rates for the problem of distributed estimation of discrete distributions under communication constraints, where $\$n\$$ users observing $\$m\$$ samples each can broadcast only $\$ell\$$ bits. Our main result is a tight characterization (up to logarithmic factors) of the error rate as a function of $\$m\$, \$ell\$,$ the domain size, and the number of users under most regimes of interest. While previous work focused on the setting where each user only holds one sample, we show that as $\$m\$$ grows the $\$ell_1\$$ error rate gets reduced by a factor of $\$sqrt{m}\$$ for small $\$m\$. However, for large $\$m\$$ we observe an interesting phase transition: the dependence of the error rate on the communication constraint $\$ell\$$ changes from $\$1/sqrt{2^{ell}}\$$ to $\$1/sqrt{ell}\$$.$

Revisiting Deep Learning Models for Tabular Data

- Yury Gorishniy, Ivan Rubachev, Valentin Khrulkov, Artem Babenko

- abstract@[open-review](#): The existing literature on deep learning for tabular data proposes a wide range of novel architectures and reports competitive results on various datasets. However, the proposed models are usually not properly compared to each other and existing works often use different benchmarks and experiment protocols. As a result, it is unclear for both researchers and practitioners what models perform best. Additionally, the field still lacks effective baselines, that is, the easy-to-use models that provide competitive performance across different problems. In this work, we perform an overview of the main families of DL architectures for tabular data and raise the bar of baselines in tabular DL by identifying two simple and powerful deep architectures. The first one is a ResNet-like architecture which turns out to be a strong baseline that is often missing in prior works. The second model is our simple adaptation of the Transformer architecture for tabular data, which outperforms other solutions on most tasks. Both models are compared to many existing architectures on a diverse set of tasks under the same training and tuning protocols. We also compare the best DL models with Gradient Boosted Decision Trees and conclude that there is still no universally superior solution. The source code is available at <https://github.com/yandex-research/rtdl>.

[Backdoor Attack with Imperceptible Input and Latent Modification](#)

- Khoa Doan, Yingjie Lao, Ping Li
- abstract@[open-review](#): Recent studies have shown that deep neural networks (DNN) are vulnerable to various adversarial attacks. In particular, an adversary can inject a stealthy backdoor into a model such that the compromised model will behave normally without the presence of the trigger. Techniques for generating backdoor images that are visually imperceptible from clean images have also been developed recently, which further enhance the stealthiness of the backdoor attacks from the input space. Along with the development of attacks, defense against backdoor attacks is also evolving. Many existing countermeasures found that backdoor tends to leave tangible footprints in the latent or feature space, which can be utilized to mitigate backdoor attacks. In this paper, we extend the concept of imperceptible backdoor from the input space to the latent representation, which significantly improves the effectiveness against the existing defense mechanisms, especially those relying on the distinguishability between clean inputs and backdoor inputs in latent space. In the proposed framework, the trigger function will learn to manipulate the input by injecting imperceptible input noise while matching the latent representations of the clean and manipulated inputs via a Wasserstein-based regularization of the corresponding empirical distributions. We formulate such an objective as a non-convex and constrained optimization problem and solve the problem with an efficient stochastic alternating optimization procedure. We name the proposed backdoor attack as Wasserstein Backdoor (WB), which achieves a high attack success rate while being stealthy from both the input and latent spaces, as tested in several benchmark datasets, including MNIST, CIFAR10, GTSRB, and TinyImagenet.

[SOPE: Spectrum of Off-Policy Estimators](#)

- Christina Yuan, Yash Chandak, Stephen Giguere, Philip S. Thomas, Scott Niekum
- abstract@[open-review](#): Many sequential decision making problems are high-stakes and require off-policy evaluation (OPE) of a new policy using historical data collected using some other policy. One of the most common OPE techniques that provides unbiased estimates is trajectory based importance sampling (IS). However, due to the high variance of trajectory IS estimates, importance sampling methods based on state-action visitation distributions (SIS) have recently been adopted. Unfortunately, while SIS often provides lower variance estimates for long horizons, estimating the state-action distribution ratios can be challenging and lead to biased estimates. In this paper, we present a new perspective on this bias-variance trade-off and show the existence of a spectrum of estimators whose endpoints are SIS and IS. Additionally, we also establish a spectrum for doubly-robust and weighted version of these estimators. We provide empirical evidence that estimators in this spectrum can be used to trade-off between the bias and variance of IS and SIS and can achieve lower mean-squared error than both IS and SIS.

[Label-Imbalanced and Group-Sensitive Classification under Overparameterization](#)

- Ganesh Ramachandra Kini, Orestis Paraskevas, Samet Oymak, Christos Thrampoulidis
- abstract@[open-review](#): The goal in label-imbalanced and group-sensitive classification is to optimize relevant metrics such as balanced error and equal opportunity. Classical methods, such as weighted cross-entropy, fail when training deep nets to the terminal phase of training (TPT), that is training beyond zero training error. This observation has motivated recent flurry of activity in developing heuristic alternatives following the intuitive mechanism of promoting larger margin for minorities. In contrast to previous heuristics, we follow a principled analysis explaining how different loss adjustments affect margins. First, we prove that for all linear classifiers trained in TPT, it is necessary to introduce multiplicative, rather than additive, logit adjustments so that the interclass margins change appropriately. To show this, we discover a connection of the multiplicative CE modification to the cost-sensitive support-vector machines. Perhaps counterintuitively, we also find that, at the start of training, the same multiplicative weights can actually harm the minority classes. Thus, while additive adjustments are ineffective in the TPT, we show that they can speed up convergence by countering the initial negative effect of the multiplicative weights. Motivated by these findings, we formulate the vector-scaling (VS) loss, that captures existing techniques as special cases. Moreover, we introduce a natural extension of the VS-loss to group-sensitive classification, thus treating the two common types of imbalances (label/group) in a unifying way. Importantly, our experiments on state-of-the-art datasets are fully consistent with our theoretical insights and confirm the superior performance of our algorithms. Finally, for imbalanced Gaussian-mixtures data, we perform a generalization analysis, revealing tradeoffs between balanced / standard error and equal opportunity.

[Neural Program Generation Modulo Static Analysis](#)

- Rohan Mukherjee, Yeming Wen, Dipak Chaudhari, Thomas Reps, Swarat Chaudhuri, Christopher Jermaine
- abstract@[open-review](#): State-of-the-art neural models of source code tend to be evaluated on the generation of individual expressions and lines of code, and commonly fail on long-horizon tasks such as the generation of entire method bodies. We propose to address this deficiency using weak supervision from a static program analyzer. Our neurosymbolic method allows a deep generative model to symbolically compute, using calls to a static analysis tool, long-distance semantic relationships in the code that it has already generated. During training, the model observes these relationships and learns to generate programs conditioned on them. We apply our approach to the problem of generating entire Java methods given the remainder of the class that contains the method. Our experiments show that the approach substantially outperforms a state-of-the-art transformer and a model that explicitly tries to learn program semantics on this task, both in terms of producing programs free of basic semantic errors and in terms of syntactically matching the ground truth.

[Unfolding Taylor's Approximations for Image Restoration](#)

- man zhou, Xueyang Fu, Zeyu Xiao, Gang Yang, Aiping Liu, Zhiwei Xiong
- abstract@[open-review](#): Deep learning provides a new avenue for image restoration, which demands a delicate balance between fine-grained details and high-level contextualized information during recovering the latent clear image. In practice, however, existing methods empirically construct encapsulated end-to-end mapping networks without deepening into the rationality, and neglect the intrinsic prior knowledge of restoration task. To solve the above problems, inspired by Taylor's Approximations, we unfold Taylor's Formula to construct a novel framework for image restoration. We find the main part and the derivative part of Taylor's Approximations take the same effect as the two competing goals of high-level contextualized information and spatial details of image restoration respectively. Specifically, our framework consists of two steps, which are correspondingly responsible for the mapping and derivative functions. The former first learns the high-level contextualized information and the later combines it with the degraded input to progressively recover local high-order spatial details. Our proposed framework is orthogonal to existing methods and thus can be easily integrated with them for further improvement, and extensive experiments demonstrate the effectiveness and scalability of our proposed framework.

[Metropolis-Hastings Data Augmentation for Graph Neural Networks](#)

- Hyeonjin Park, Seunghun Lee, Sihyeon Kim, Jinyoung Park, Jisu Jeong, Kyung-Min Kim, Jung-Woo Ha, Hyunwoo J. Kim
- abstract@[open-review](#): Graph Neural Networks (GNNs) often suffer from weak-generalization due to sparsely labeled data despite their promising results on various graph-based tasks. Data augmentation is a prevalent remedy to improve the generalization ability of models in many domains. However, due to the non-Euclidean nature of data space and the dependencies between samples, designing effective augmentation on graphs is challenging. In this paper, we propose a novel framework Metropolis-Hastings Data Augmentation (MH-Aug) that draws augmented graphs from an explicit target distribution for semi-supervised learning. MH-Aug produces a sequence of augmented graphs from the target distribution enables flexible control of the strength and diversity of augmentation. Since the direct sampling from the complex target distribution is challenging, we adopt the Metropolis-Hastings algorithm to obtain the augmented samples. We also propose a simple and effective semi-supervised learning strategy with generated samples from MH-Aug. Our extensive experiments demonstrate that MH-Aug can generate a sequence of samples according to the target distribution to significantly improve the performance of GNNs.

[Strategic Behavior is Bliss: Iterative Voting Improves Social Welfare](#)

- Joshua Kavner, Lirong Xia
- abstract@[open-review](#): Recent work in iterative voting has defined the additive dynamic price of anarchy (ADPoA) as the difference in social welfare between the truthful and worst-case equilibrium profiles resulting from repeated strategic manipulations. While iterative plurality has been shown to only return alternatives with at most one less initial votes than the truthful winner, it is less understood how agents' welfare changes in equilibrium. To this end, we differentiate agents' utility from their manipulation mechanism and determine iterative plurality's ADPoA in the worst- and average-cases. We first prove that the worst-case ADPoA is linear in the number of agents. To overcome this negative result, we study the average-case ADPoA and prove that equilibrium winners have a constant order welfare advantage over the truthful winner in expectation. Our positive results illustrate the prospect for social welfare to increase due to strategic manipulation.

[Agnostic Reinforcement Learning with Low-Rank MDPs and Rich Observations](#)

- Ayush Sekhari, Christoph Dann, Mehryar Mohri, Yishay Mansour, Karthik Sridharan
- abstract@[open-review](#): There have been many recent advances on provably efficient Reinforcement Learning (RL) in problems with rich observation spaces. However, all these works share a strong realizability assumption about the optimal value function of the true MDP. Such realizability assumptions are often too strong to hold in practice. In this work, we consider the more realistic setting of agnostic RL with rich observation spaces and a fixed class of policies $\$P_i\$$ that may not contain any near-optimal policy. We provide an algorithm for this setting whose error is bounded in terms of the rank $\$d\$$ of the underlying MDP. Specifically, our algorithm enjoys a sample complexity bound of $\$widetilde{O}((H^{4d} K^{3d} \log |P|)^{\epsilon^2})\$$ where $\$H\$$ is the length of episodes, $\$K\$$ is the number of actions and $\$epsilon > 0\$$ is the desired sub-optimality. We also provide a nearly matching lower bound for this agnostic setting that shows that the exponential dependence on rank is unavoidable, without further assumptions.

[Functional Regularization for Reinforcement Learning via Learned Fourier Features](#)

- Alexander Li, Deepak Pathak
- abstract@[open-review](#): We propose a simple architecture for deep reinforcement learning by embedding inputs into a learned Fourier basis and show that it improves the sample efficiency of both state-based and image-based RL. We perform infinite-width analysis of our architecture using the Neural Tangent Kernel and theoretically show that tuning the initial variance of the Fourier basis is equivalent to functional regularization of the learned deep network. That is, these learned Fourier features allow for adjusting the degree to which networks underfit or overfit different frequencies in the training data, and hence provide a controlled mechanism to improve the stability and performance of RL optimization. Empirically, this allows us to prioritize learning low-frequency functions and speed up learning by reducing networks' susceptibility to noise in the optimization process, such as during Bellman updates. Experiments on standard state-based and image-based RL benchmarks show clear benefits of our architecture over the baselines.

[Adaptive First-Order Methods Revisited: Convex Minimization without Lipschitz Requirements](#)

- Kimon Antonakopoulos, Panayotis Mertikopoulos
- abstract@[open-review](#): We propose a new family of adaptive first-order methods for a class of convex minimization problems that may fail to be Lipschitz continuous or smooth in the standard sense. Specifically, motivated by a recent flurry of activity on non-Lipschitz (NoLips) optimization, we consider problems that are continuous or smooth relative to a reference Bregman function – as opposed to a global, ambient norm (Euclidean or otherwise). These conditions encompass a wide range of problems with singular objective, such as Fisher markets, Poisson tomography, D-design, and the like. In this setting, the application of existing order-optimal adaptive methods – like UnixGrad or AcceleGrad – is not possible, especially in the presence of randomness and uncertainty. The proposed method, adaptive mirror descent (AdaMir), aims to close this gap by concurrently achieving min-max optimal rates in problems that are relatively continuous or smooth, including stochastic ones.

[Adapting to function difficulty and growth conditions in private optimization](#)

- Hilal Asi, Daniel Levy, John C. Duchi
- abstract@[open-review](#): We develop algorithms for private stochastic convex optimization that adapt to the hardness of the specific function we wish to optimize. While previous work provide worst-case bounds for arbitrary convex functions, it is often the case that the function at hand belongs to a smaller class that enjoys faster rates. Concretely, we show that for functions exhibiting $\$kappa$-growth around the optimum, i.e., $\$f(x) \geq f(x^*) + \lambda kappa^{-1} \|x - x^*\|^2 kappa$$ for $\$kappa > 1$$, our algorithms improve upon the standard $\$sqrt{d}/(n\epsilon)$$ privacy rate to the faster $\$(sqrt{d}/(n\epsilon))^{\frac{kappa}{kappa-1}}$. Crucially, they achieve these rates without knowledge of the growth constant $\$kappa$$ of the function. Our algorithms build upon the inverse sensitivity mechanism, which adapts to instance difficulty [2], and recent localization techniques in private optimization [25]. We complement our algorithms with matching lower bounds for these function classes and demonstrate that our adaptive algorithm is simultaneously (minimax) optimal over all $\$kappa \geq 1+c$$ whenever $\$c = \Theta(1)$$.$

[Support Recovery of Sparse Signals from a Mixture of Linear Measurements](#)

- Soumyabrata Pal, Arya Mazumdar, Venkata Gandikota
- abstract@[open-review](#): Recovery of support of a sparse vector from simple measurements is a widely studied problem, considered under the frameworks of compressed sensing, 1-bit compressed sensing, and more general single index models. We consider generalizations of this problem: mixtures of linear regressions, and mixtures of linear classifiers, where the goal is to recover supports of multiple sparse vectors using only a small number of possibly noisy linear, and 1-bit measurements respectively. The key challenge is that the measurements from different vectors are randomly mixed. Both of these problems have also received attention recently. In mixtures of linear classifiers, an observation corresponds to the side of the queried hyperplane a random unknown vector lies in; whereas in mixtures of linear regressions we observe the projection of a random unknown vector on the queried hyperplane. The primary step in recovering the unknown vectors from the mixture is to first identify the support of all the individual component vectors. In this work, we study the number of measurements sufficient for recovering the supports of all the component vectors in a mixture in both these models. We provide

algorithms that use a number of measurements polynomial in $k, \log n$ and quasi-polynomial in ℓ , to recover the support of all the ℓ unknown vectors in the mixture with high probability when each individual component is a k -sparse n -dimensional vector.

[Stochastic Gradient Descent-Ascent and Consensus Optimization for Smooth Games: Convergence Analysis under Expected Co-coercivity](#)

- Nicolas Loizou, Hugo Berard, Gauthier Gidel, Ioannis Mitliagkas, Simon Lacoste-Julien
- abstract@[open-review](#): Two of the most prominent algorithms for solving unconstrained smooth games are the classical stochastic gradient descent-ascent (SGDA) and the recently introduced stochastic consensus optimization (SCO) [Mescheder et al., 2017]. SGDA is known to converge to a stationary point for specific classes of games, but current convergence analyses require a bounded variance assumption. SCO is used successfully for solving large-scale adversarial problems, but its convergence guarantees are limited to its deterministic variant. In this work, we introduce the expected co-coercivity condition, explain its benefits, and provide the first last-iterate convergence guarantees of SGDA and SCO under this condition for solving a class of stochastic variational inequality problems that are potentially non-monotone. We prove linear convergence of both methods to a neighborhood of the solution when they use constant step-size, and we propose insightful stepsize-switching rules to guarantee convergence to the exact solution. In addition, our convergence guarantees hold under the arbitrary sampling paradigm, and as such, we give insights into the complexity of minibatching.

[Tighter Expected Generalization Error Bounds via Wasserstein Distance](#)

- Borja Rodríguez Gálvez, German Bassi, Ragnar Thobaben, Mikael Skoglund
- abstract@[open-review](#): This work presents several expected generalization error bounds based on the Wasserstein distance. More specifically, it introduces full-dataset, single-letter, and random-subset bounds, and their analogous in the randomized subsample setting from Steinke and Zakynthinou [1]. Moreover, when the loss function is bounded and the geometry of the space is ignored by the choice of the metric in the Wasserstein distance, these bounds recover from below (and thus, are tighter than) current bounds based on the relative entropy. In particular, they generate new, non-vacuous bounds based on the relative entropy. Therefore, these results can be seen as a bridge between works that account for the geometry of the hypothesis space and those based on the relative entropy, which is agnostic to such geometry. Furthermore, it is shown how to produce various new bounds based on different information measures (e.g., the lautm information or several f -divergences) based on these bounds and how to derive similar bounds with respect to the backward channel using the presented proof techniques.

[Unifying Width-Reduced Methods for Quasi-Self-Concordant Optimization](#)

- Deeksha Adil, Brian Bullins, Sushant Sachdeva
- abstract@[open-review](#): We provide several algorithms for constrained optimization of a large class of convex problems, including softmax, ℓ_p regression, and logistic regression. Central to our approach is the notion of width reduction, a technique which has proven immensely useful in the context of maximum flow [Christiano et al., STOC'11] and, more recently, ℓ_p regression [Adil et al., SODA'19], in terms of improving the iteration complexity from $O(m^{1/2})$ to $\tilde{O}(m^{1/3})$, where m is the number of rows of the design matrix, and where each iteration amounts to a linear system solve. However, a considerable drawback is that these methods require both problem-specific potentials and individually tailored analyses. As our main contribution, we initiate a new direction of study by presenting the first unified approach to achieving $m^{1/3}$ -type rates. Notably, our method goes beyond these previously considered problems to more broadly capture quasi-self-concordant losses, a class which has recently generated much interest and includes the well-studied problem of logistic regression, among others. In order to do so, we develop a unified width reduction method for carefully handling these losses based on a more general set of potentials. Additionally, we directly achieve $m^{1/3}$ -type rates in the constrained setting without the need for any explicit acceleration schemes, thus naturally complementing recent work based on a ball-oracle approach [Carmon et al., NeurIPS'20].

[Bridging the Imitation Gap by Adaptive Insubordination](#)

- Luca Weihs, Unnat Jain, Iou-Jen Liu, Jordi Salvador, Svetlana Lazebnik, Aniruddha Kembhavi, Alex Schwing
- abstract@[open-review](#): In practice, imitation learning is preferred over pure reinforcement learning whenever it is possible to design a teaching agent to provide expert supervision. However, we show that when the teaching agent makes decisions with access to privileged information that is unavailable to the student, this information is marginalized during imitation learning, resulting in an "imitation gap" and, potentially, poor results. Prior work bridges this gap via a progression from imitation learning to reinforcement learning. While often successful, gradual progression fails for tasks that require frequent switches between exploration and memorization. To better address these tasks and alleviate the imitation gap we propose 'Adaptive Insubordination' (ADVISOR). ADVISOR dynamically weights imitation and reward-based reinforcement learning losses during training, enabling on-the-fly switching between imitation and exploration. On a suite of challenging tasks set within gridworlds, multi-agent particle environments, and high-fidelity 3D simulators, we show that on-the-fly switching with ADVISOR outperforms pure imitation, pure reinforcement learning, as well as their sequential and parallel combinations.

[Adversarial Robustness with Non-uniform Perturbations](#)

- Ecenaz Erdemir, Jeffrey Bickford, Luca Melis, Sergul Aydore
- abstract@[open-review](#): Robustness of machine learning models is critical for security related applications, where real-world adversaries are uniquely focused on evading neural network based detectors. Prior work mainly focus on crafting adversarial examples (AEs) with small uniform norm-bounded perturbations across features to maintain the requirement of imperceptibility. However, uniform perturbations do not result in realistic AEs in domains such as malware, finance, and social networks. For these types of applications, features typically have some semantically meaningful dependencies. The key idea of our proposed approach is to enable non-uniform perturbations that can adequately represent these feature dependencies during adversarial training. We propose using characteristics of the empirical data distribution, both on correlations between the features and the importance of the features themselves. Using experimental datasets for malware classification, credit risk prediction, and spam detection, we show that our approach is more robust to real-world attacks. Finally, we present robustness certification utilizing non-uniform perturbation bounds, and show that non-uniform bounds achieve better certification.

[Container: Context Aggregation Networks](#)

- peng gao, Jiasen Lu, hongsheng Li, Roozbeh Mottaghi, Aniruddha Kembhavi
- abstract@[open-review](#): Convolutional neural networks (CNNs) are ubiquitous in computer vision, with a myriad of effective and efficient variations. Recently, Transformers -- originally introduced in natural language processing -- have been increasingly adopted in computer vision. While early adopters continued to employ CNN backbones, the latest networks are end-to-end CNN-free Transformer solutions. A recent surprising finding now shows that a simple MLP based solution without any traditional convolutional or Transformer components can produce effective visual representations. While CNNs, Transformers and MLP-Mixers may be considered as completely disparate architectures, we provide a unified view showing that they are in fact special cases of a more general method to aggregate spatial context in a neural network stack. We present the model (CONText Aggregation NEtwork), a general-purpose building block for multi-head context aggregation that can exploit long-range interactions a la Transformers while still exploiting the inductive bias of the local convolution operation leading to faster convergence speeds, often seen in CNNs. Our model architecture achieves 82.7% Top-1 accuracy on ImageNet using 22M parameters, +2.8 improvement compared with DeiT-Small, and can converge to 79.9% Top-1 accuracy in just 200 epochs. In contrast to Transformer-based methods that do not scale well to downstream tasks that rely on larger input image resolutions, our efficient

network, named \modellight, can be employed in object detection and instance segmentation networks such as DETR, RetinaNet and Mask-RCNN to obtain an impressive detection mAP of 38.9, 43.8, 45.1 and mask mAP of 41.3, providing large improvements of 6.6, 7.3, 6.9 and 6.6 pts respectively, compared to a ResNet-50 backbone with a comparable compute and parameter size. Our method also achieves promising results on self-supervised learning compared to DeiT on the DINO framework. Code is released at <https://github.com/allenai/container>.

[ConE: Cone Embeddings for Multi-Hop Reasoning over Knowledge Graphs](#)

- Zhanqiu Zhang, Jie Wang, Jiajun Chen, Shuiwang Ji, Feng Wu
- abstract@[open-review](#): Query embedding (QE)---which aims to embed entities and first-order logical (FOL) queries in low-dimensional spaces---has shown great power in multi-hop reasoning over knowledge graphs. Recently, embedding entities and queries with geometric shapes becomes a promising direction, as geometric shapes can naturally represent answer sets of queries and logical relationships among them. However, existing geometry-based models have difficulty in modeling queries with negation, which significantly limits their applicability. To address this challenge, we propose a novel query embedding model, namely \textbf{Con}e \textbf{E}mbbeddings (ConE), which is the first geometry-based QE model that can handle all the FOL operations, including conjunction, disjunction, and negation. Specifically, ConE represents entities and queries as Cartesian products of two-dimensional cones, where the intersection and union of cones naturally model the conjunction and disjunction operations. By further noticing that the closure of complement of cones remains cones, we design geometric complement operators in the embedding space for the negation operations. Experiments demonstrate that ConE significantly outperforms existing state-of-the-art methods on benchmark datasets.

[Federated Hyperparameter Tuning: Challenges, Baselines, and Connections to Weight-Sharing](#)

- Mikhail Khodak, Renbo Tu, Tian Li, Liam Li, Maria-Florina F. Balcan, Virginia Smith, Ameet Talwalkar
- abstract@[open-review](#): Tuning hyperparameters is a crucial but arduous part of the machine learning pipeline. Hyperparameter optimization is even more challenging in federated learning, where models are learned over a distributed network of heterogeneous devices; here, the need to keep data on device and perform local training makes it difficult to efficiently train and evaluate configurations. In this work, we investigate the problem of federated hyperparameter tuning. We first identify key challenges and show how standard approaches may be adapted to form baselines for the federated setting. Then, by making a novel connection to the neural architecture search technique of weight-sharing, we introduce a new method, FedEx, to accelerate federated hyperparameter tuning that is applicable to widely-used federated optimization methods such as FedAvg and recent variants. Theoretically, we show that a FedEx variant correctly tunes the on-device learning rate in the setting of online convex optimization across devices. Empirically, we show that FedEx can outperform natural baselines for federated hyperparameter tuning by several percentage points on the Shakespeare, FEMNIST, and CIFAR-10 benchmarks---obtaining higher accuracy using the same training budget.

[Training for the Future: A Simple Gradient Interpolation Loss to Generalize Along Time](#)

- Anshul Nasery, Soumyadeep Thakur, Vihari Piratla, Abir De, Sunita Sarawagi
- abstract@[open-review](#): In several real world applications, machine learning models are deployed to make predictions on data whose distribution changes gradually along time, leading to a drift between the train and test distributions. Such models are often re-trained on new data periodically, and they hence need to generalize to data not too far into the future. In this context, there is much prior work on enhancing temporal generalization, e.g. continuous transportation of past data, kernel smoothed time-sensitive parameters and more recently, adversarial learning of time-invariant features. However, these methods share several limitations, e.g., poor scalability, training instability, and dependence on unlabeled data from the future. Responding to the above limitations, we propose a simple method that starts with a model with time-sensitive parameters but regularizes its temporal complexity using a Gradient Interpolation (GI) loss. GI allows the decision boundary to change along time and can still prevent overfitting to the limited training time snapshots by allowing task-specific control over changes along time. We compare our method to existing baselines on multiple real-world datasets, which show that GI outperforms more complicated generative and adversarial approaches on the one hand, and simpler gradient regularization methods on the other.

[Agent Modelling under Partial Observability for Deep Reinforcement Learning](#)

- Georgios Papoudakis, Filippos Christianos, Stefano Albrecht
- abstract@[open-review](#): Modelling the behaviours of other agents is essential for understanding how agents interact and making effective decisions. Existing methods for agent modelling commonly assume knowledge of the local observations and chosen actions of the modelled agents during execution. To eliminate this assumption, we extract representations from the local information of the controlled agent using encoder-decoder architectures. Using the observations and actions of the modelled agents during training, our models learn to extract representations about the modelled agents conditioned only on the local observations of the controlled agent. The representations are used to augment the controlled agent's decision policy which is trained via deep reinforcement learning; thus, during execution, the policy does not require access to other agents' information. We provide a comprehensive evaluation and ablations studies in cooperative, competitive and mixed multi-agent environments, showing that our method achieves significantly higher returns than baseline methods which do not use the learned representations.

[Leveraging Distribution Alignment via Stein Path for Cross-Domain Cold-Start Recommendation](#)

- Weiming Liu, Jiajie Su, Chaochao Chen, Xiaolin Zheng
- abstract@[open-review](#): Cross-Domain Recommendation (CDR) has been popularly studied to utilize different domain knowledge to solve the cold-start problem in recommender systems. In this paper, we focus on the Cross-Domain Cold-Start Recommendation (CDCSR) problem. That is, how to leverage the information from a source domain, where items are 'warm', to improve the recommendation performance of a target domain, where items are 'cold'. Unfortunately, previous approaches on cold-start and CDR cannot reduce the latent embedding discrepancy across domains efficiently and lead to model degradation. To address this issue, we propose DisAlign, a cross-domain recommendation framework for the CDCSR problem, which utilizes both rating and auxiliary representations from the source domain to improve the recommendation performance of the target domain. Specifically, we first propose Stein path alignment for aligning the latent embedding distributions across domains, and then further propose its improved version, i.e., proxy Stein path, which can reduce the operation consumption and improve efficiency. Our empirical study on Douban and Amazon datasets demonstrate that DisAlign significantly outperforms the state-of-the-art models under the CDCSR setting.

[Conservative Offline Distributional Reinforcement Learning](#)

- Yecheng Ma, Dinesh Jayaraman, Osbert Bastani
- abstract@[open-review](#): Many reinforcement learning (RL) problems in practice are offline, learning purely from observational data. A key challenge is how to ensure the learned policy is safe, which requires quantifying the risk associated with different actions. In the online setting, distributional RL algorithms do so by learning the distribution over returns (i.e., cumulative rewards) instead of the expected return; beyond quantifying risk, they have also been shown to learn better representations for planning. We propose Conservative Offline Distributional Actor Critic (CODAC), an offline RL algorithm suitable for both risk-neutral and risk-averse domains. CODAC adapts distributional RL to the offline setting by penalizing the predicted quantiles of the return for out-of-distribution actions. We prove that CODAC learns a conservative return distribution---in particular, for finite MDPs, CODAC converges to an uniform lower bound on the quantiles of the return distribution; our proof relies on a novel analysis of the distributional Bellman operator. In our experiments, on two challenging robot navigation tasks, CODAC successfully learns risk-averse policies using offline data collected purely from risk-neutral agents. Furthermore, CODAC is state-of-the-art on the D4RL MuJoCo benchmark in terms of both expected and risk-sensitive performance.

[Separation Results between Fixed-Kernel and Feature-Learning Probability Metrics](#)

- Carles Domingo i Enrich, Youssef Mroueh
- abstract@[open-review](#): Several works in implicit and explicit generative modeling empirically observed that feature-learning discriminators outperform fixed-kernel discriminators in terms of the sample quality of the models. We provide separation results between probability metrics with fixed-kernel and feature-learning discriminators using the function classes \mathcal{F}_2 and \mathcal{F}_1 respectively, which were developed to study overparametrized two-layer neural networks. In particular, we construct pairs of distributions over hyper-spheres that can not be discriminated by fixed kernel \mathcal{F}_2 integral probability metric (IPM) and Stein discrepancy (SD) in high dimensions, but that can be discriminated by their feature learning (\mathcal{F}_1) counterparts. To further study the separation we provide links between the \mathcal{F}_1 and \mathcal{F}_2 IPMs with sliced Wasserstein distances. Our work suggests that fixed-kernel discriminators perform worse than their feature learning counterparts because their corresponding metrics are weaker.

[Risk Minimization from Adaptively Collected Data: Guarantees for Supervised and Policy Learning](#)

- Aurelien Bibaut, Nathan Kallus, Maria Dimakopoulou, Antoine Chambaz, Mark van der Laan
- abstract@[open-review](#): Empirical risk minimization (ERM) is the workhorse of machine learning, whether for classification and regression or for off-policy policy learning, but its model-agnostic guarantees can fail when we use adaptively collected data, such as the result of running a contextual bandit algorithm. We study a generic importance sampling weighted ERM algorithm for using adaptively collected data to minimize the average of a loss function over a hypothesis class and provide first-of-their-kind generalization guarantees and fast convergence rates. Our results are based on a new maximal inequality that carefully leverages the importance sampling structure to obtain rates with the good dependence on the exploration rate in the data. For regression, we provide fast rates that leverage the strong convexity of squared-error loss. For policy learning, we provide regret guarantees that close an open gap in the existing literature whenever exploration decays to zero, as is the case for bandit-collected data. An empirical investigation validates our theory.

[Bayesian Optimization with High-Dimensional Outputs](#)

- Wesley J. Maddox, Maximilian Balandat, Andrew G. Wilson, Eytan Bakshy
- abstract@[open-review](#): Bayesian optimization is a sample-efficient black-box optimization procedure that is typically applied to a small number of independent objectives. However, in practice we often wish to optimize objectives defined over many correlated outcomes (or “tasks”). For example, scientists may want to optimize the coverage of a cell tower network across a dense grid of locations. Similarly, engineers may seek to balance the performance of a robot across dozens of different environments via constrained or robust optimization. However, the Gaussian Process (GP) models typically used as probabilistic surrogates for multi-task Bayesian optimization scale poorly with the number of outcomes, greatly limiting applicability. We devise an efficient technique for exact multi-task GP sampling that combines exploiting Kronecker structure in the covariance matrices with Matheron’s identity, allowing us to perform Bayesian optimization using exact multi-task GP models with tens of thousands of correlated outputs. In doing so, we achieve substantial improvements in sample efficiency compared to existing approaches that model solely the outcome metrics. We demonstrate how this unlocks a new class of applications for Bayesian optimization across a range of tasks in science and engineering, including optimizing interference patterns of an optical interferometer with 65,000 outputs.

[Finding Optimal Tangent Points for Reducing Distortions of Hard-label Attacks](#)

- Chen Ma, Xiangyu Guo, Li Chen, Jun-Hai Yong, Yisen Wang
- abstract@[open-review](#): One major problem in black-box adversarial attacks is the high query complexity in the hard-label attack setting, where only the top-1 predicted label is available. In this paper, we propose a novel geometric-based approach called Tangent Attack (TA), which identifies an optimal tangent point of a virtual hemisphere located on the decision boundary to reduce the distortion of the attack. Assuming the decision boundary is locally flat, we theoretically prove that the minimum ℓ_2 distortion can be obtained by reaching the decision boundary along the tangent line passing through such tangent point in each iteration. To improve the robustness of our method, we further propose a generalized method which replaces the hemisphere with a semi-ellipsoid to adapt to curved decision boundaries. Our approach is free of pre-training. Extensive experiments conducted on the ImageNet and CIFAR-10 datasets demonstrate that our approach can consume only a small number of queries to achieve the low-magnitude distortion. The implementation source code is released online.

[Scalable Diverse Model Selection for Accessible Transfer Learning](#)

- Daniel Bolya, Rohit Mittapalli, Judy Hoffman
- abstract@[open-review](#): With the preponderance of pretrained deep learning models available off-the-shelf from model banks today, finding the best weights to fine-tune to your use-case can be a daunting task. Several methods have recently been proposed to find good models for transfer learning, but they either don’t scale well to large model banks or don’t perform well on the diversity of off-the-shelf models. Ideally the question we want to answer is, “given some data and a source model, can you quickly predict the model’s accuracy after fine-tuning?” In this paper, we formalize this setting as “Scalable Diverse Model Selection” and propose several benchmarks for evaluating on this task. We find that existing model selection and transferability estimation methods perform poorly here and analyze why this is the case. We then introduce simple techniques to improve the performance and speed of these algorithms. Finally, we iterate on existing methods to create PARC, which outperforms all other methods on diverse model selection. We have released the benchmarks and method code in hope to inspire future work in model selection for accessible transfer learning.

[Light Field Networks: Neural Scene Representations with Single-Evaluation Rendering](#)

- Vincent Sitzmann, Semon Rezchikov, Bill Freeman, Josh Tenenbaum, Fredo Durand
- abstract@[open-review](#): Inferring representations of 3D scenes from 2D observations is a fundamental problem of computer graphics, computer vision, and artificial intelligence. Emerging 3D-structured neural scene representations are a promising approach to 3D scene understanding. In this work, we propose a novel neural scene representation, Light Field Networks or LFNs, which represent both geometry and appearance of the underlying 3D scene in a 360-degree, four-dimensional light field parameterized via a neural implicit representation. Rendering a ray from an LFN requires only a single network evaluation, as opposed to hundreds of evaluations per ray for ray-marching or volumetric based renderers in 3D-structured neural scene representations. In the setting of simple scenes, we leverage meta-learning to learn a prior over LFNs that enables multi-view consistent light field reconstruction from as little as a single image observation. This results in dramatic reductions in time and memory complexity, and enables real-time rendering. The cost of storing a 360-degree light field via an LFN is two orders of magnitude lower than conventional methods such as the Lumigraph. Utilizing the analytical differentiability of neural implicit representations and a novel parameterization of light space, we further demonstrate the extraction of sparse depth maps from LFNs.

[ViSER: Video-Specific Surface Embeddings for Articulated 3D Shape Reconstruction](#)

- Gengshan Yang, Deqing Sun, Varun Jampani, Daniel Vlasic, Forrester Cole, Ce Liu, Deva Ramanan
- abstract@[open-review](#): We introduce ViSER, a method for recovering articulated 3D shapes and dense 3D trajectories from monocular videos. Previous work on high-quality reconstruction of dynamic 3D shapes typically relies on multiple camera views, strong category-specific priors, or 2D keypoint

supervision. We show that none of these are required if one can reliably estimate long-range correspondences in a video, making use of only 2D object masks and two-frame optical flow as inputs. ViSER infers correspondences by matching 2D pixels to a canonical, deformable 3D mesh via video-specific surface embeddings that capture the pixel appearance of each surface point. These embeddings behave as a continuous set of keypoint descriptors defined over the mesh surface, which can be used to establish dense long-range correspondences across pixels. The surface embeddings are implemented as coordinate-based MLPs that are fit to each video via self-supervised losses. Experimental results show that ViSER compares favorably against prior work on challenging videos of humans with loose clothing and unusual poses as well as animals videos from DAVIS and YTVOS. Project page: viser-shape.github.io.

[Understanding the Effect of Stochasticity in Policy Optimization](#)

- Jincheng Mei, Bo Dai, Chenjun Xiao, Csaba Szepesvari, Dale Schuurmans
- abstract@[open-review](#): We study the effect of stochasticity in on-policy policy optimization, and make the following four contributions. \emph{First}, we show that the preferability of optimization methods depends critically on whether stochastic versus exact gradients are used. In particular, unlike the true gradient setting, geometric information cannot be easily exploited in the stochastic case for accelerating policy optimization without detrimental consequences or impractical assumptions. \emph{Second}, to explain these findings we introduce the concept of committal rate for stochastic policy optimization, and show that this can serve as a criterion for determining almost sure convergence to global optimality. \emph{Third}, we show that in the absence of external oracle information, which allows an algorithm to determine the difference between optimal and sub-optimal actions given only on-policy samples, there is an inherent trade-off between exploiting geometry to accelerate convergence versus achieving optimality almost surely. That is, an uninformed algorithm either converges to a globally optimal policy with probability \$1\$ but at a rate no better than \$O(1/t)\$, or it achieves faster than \$O(1/t)\$ convergence but then must fail to converge to the globally optimal policy with some positive probability. \emph{Finally}, we use the committal rate theory to explain why practical policy optimization methods are sensitive to random initialization, then develop an ensemble method that can be guaranteed to achieve near-optimal solutions with high probability.

[Fine-Grained Zero-Shot Learning with DNA as Side Information](#)

- Sarkhan Badirli, Zeynep Akata, George Mohler, Christine Picard, Mehmet M Dundar
- abstract@[open-review](#): Fine-grained zero-shot learning task requires some form of side-information to transfer discriminative information from seen to unseen classes. As manually annotated visual attributes are extremely costly and often impractical to obtain for a large number of classes, in this study we use DNA as a side information for the first time for fine-grained zero-shot classification of species. Mitochondrial DNA plays an important role as a genetic marker in evolutionary biology and has been used to achieve near perfect accuracy in species classification of living organisms. We implement a simple hierarchical Bayesian model that uses DNA information to establish the hierarchy in the image space and employs local priors to define surrogate classes for unseen ones. On the benchmark CUB dataset we show that DNA can be equally promising, yet in general a more accessible alternative than word vectors as a side information. This is especially important as obtaining robust word representations for fine-grained species names is not a practicable goal when information about these species in free-form text is limited. On a newly compiled fine-grained insect dataset that uses DNA information from over a thousand species we show that the Bayesian approach outperforms state-of-the-art by a wide margin.

[Optimal Underdamped Langevin MCMC Method](#)

- Zhengmian Hu, Feihu Huang, Heng Huang
- abstract@[open-review](#): In the paper, we study the underdamped Langevin diffusion (ULD) with strongly-convex potential consisting of finite summation of \$N\$ smooth components, and propose an efficient discretization method, which requires \$O(N+d^{1/3}N^{2/3}\varepsilon^{2/3})\$ gradient evaluations to achieve \$\varepsilon\$-error (in \$\sqrt{\mathbb{E}\|\cdot\|_2^2}\$ distance) for approximating \$d\$-dimensional ULD. Moreover, we prove a lower bound of gradient complexity as \$\Omega(\Omega(N+d^{1/3}N^{2/3}\varepsilon^{2/3}))\$, which indicates that our method is optimal in dependence of \$N\$, \$\varepsilon\$, and \$d\$. In particular, we apply our method to sample the strongly-log-concave distribution and obtain gradient complexity better than all existing gradient based sampling algorithms. Experimental results on both synthetic and real-world data show that our new method consistently outperforms the existing ULD approaches.

[Scheduling jobs with stochastic holding costs](#)

- Dabeen Lee, Milan Vojnovic
- abstract@[open-review](#): This paper proposes a learning and scheduling algorithm to minimize the expected cumulative holding cost incurred by jobs, where statistical parameters defining their individual holding costs are unknown a priori. In each time slot, the server can process a job while receiving the realized random holding costs of the jobs remaining in the system. Our algorithm is a learning-based variant of the \$\text{cmu}\$ rule for scheduling: it starts with a preemption period of fixed length which serves as a learning phase, and after accumulating enough data about individual jobs, it switches to nonpreemptive scheduling mode. The algorithm is designed to handle instances with large or small gaps in jobs' parameters and achieves near-optimal performance guarantees. The performance of our algorithm is captured by its regret, where the benchmark is the minimum possible cost attained when the statistical parameters of jobs are fully known. We prove upper bounds on the regret of our algorithm, and we derive a regret lower bound that is almost matching the proposed upper bounds. Our numerical results demonstrate the effectiveness of our algorithm and show that our theoretical regret analysis is nearly tight.

[REMIPS: Physically Consistent 3D Reconstruction of Multiple Interacting People under Weak Supervision](#)

- Mihai Fieraru, Mihai Zanfir, Teodor Szente, Eduard Bazavan, Vlad Olaru, Cristian Sminchisescu
- abstract@[open-review](#): The three-dimensional reconstruction of multiple interacting humans given a monocular image is crucial for the general task of scene understanding, as capturing the subtleties of interaction is often the very reason for taking a picture. Current 3D human reconstruction methods either treat each person independently, ignoring most of the context, or reconstruct people jointly, but cannot recover interactions correctly when people are in close proximity. In this work, we introduce \textbf{REMIPS}, a model for 3D \underline{Re}construction of \underline{M}ultiple \underline{I}nteracting \underline{P}eople under Weak \underline{S}upervision. \textbf{REMIPS} can reconstruct a variable number of people directly from monocular images. At the core of our methodology stands a novel transformer network that combines unordered person tokens (one for each detected human) with positional-encoded tokens from image features patches. We introduce a novel unified model for self- and interpenetration-collisions based on a mesh approximation computed by applying decimation operators. We rely on self-supervised losses for flexibility and generalisation in-the-wild and incorporate self-contact and interaction-contact losses directly into the learning process. With \textbf{REMIPS}, we report state-of-the-art quantitative results on common benchmarks even in cases where no 3D supervision is used. Additionally, qualitative visual results show that our reconstructions are plausible in terms of pose and shape and coherent for challenging images, collected in-the-wild, where people are often interacting.

[Differentiable Annealed Importance Sampling and the Perils of Gradient Noise](#)

- Guodong Zhang, Kyle Hsu, Jianing Li, Chelsea Finn, Roger B. Grosse
- abstract@[open-review](#): Annealed importance sampling (AIS) and related algorithms are highly effective tools for marginal likelihood estimation, but are not fully differentiable due to the use of Metropolis-Hastings correction steps. Differentiability is a desirable property as it would admit the possibility of optimizing marginal likelihood as an objective using gradient-based methods. To this end, we propose Differentiable AIS (DAIS), a variant of AIS which ensures differentiability by abandoning the Metropolis-Hastings corrections. As a further advantage, DAIS allows for mini-batch gradients. We provide a

detailed convergence analysis for Bayesian linear regression which goes beyond previous analyses by explicitly accounting for the sampler not having reached equilibrium. Using this analysis, we prove that DAIS is consistent in the full-batch setting and provide a sublinear convergence rate. Furthermore, motivated by the problem of learning from large-scale datasets, we study a stochastic variant of DAIS that uses mini-batch gradients. Surprisingly, stochastic DAIS can be arbitrarily bad due to a fundamental incompatibility between the goals of last-iterate convergence to the posterior and elimination of the accumulated stochastic error. This is in stark contrast with other settings such as gradient-based optimization and Langevin dynamics, where the effect of gradient noise can be washed out by taking smaller steps. This indicates that annealing-based marginal likelihood estimation with stochastic gradients may require new ideas.

[PSD Representations for Effective Probability Models](#)

- Alessandro Rudi, Carlo Ciliberto
- abstract@[open-review](#): Finding a good way to model probability densities is key to probabilistic inference. An ideal model should be able to concisely approximate any probability while being also compatible with two main operations: multiplications of two models (product rule) and marginalization with respect to a subset of the random variables (sum rule). In this work, we show that a recently proposed class of positive semi-definite (PSD) models for non-negative functions is particularly suited to this end. In particular, we characterize both approximation and generalization capabilities of PSD models, showing that they enjoy strong theoretical guarantees. Moreover, we show that we can perform efficiently both sum and product rule in closed form via matrix operations, enjoying the same versatility of mixture models. Our results open the way to applications of PSD models to density estimation, decision theory, and inference.

[Exploiting a Zoo of Checkpoints for Unseen Tasks](#)

- Jiaji Huang, Qiang Qiu, Kenneth Church
- abstract@[open-review](#): There are so many models in the literature that it is difficult for practitioners to decide which combinations are likely to be effective for a new task. This paper attempts to address this question by capturing relationships among checkpoints published on the web. We model the space of tasks as a Gaussian process. The covariance can be estimated from checkpoints and unlabeled probing data. With the Gaussian process, we can identify representative checkpoints by a maximum mutual information criterion. This objective is submodular. A greedy method identifies representatives that are likely to "cover" the task space. These representatives generalize to new tasks with superior performance. Empirical evidence is provided for applications from both computational linguistics as well as computer vision.

[Towards Open-World Feature Extrapolation: An Inductive Graph Learning Approach](#)

- Qitian Wu, Chenxiao Yang, Junchi Yan
- abstract@[open-review](#): We target open-world feature extrapolation problem where the feature space of input data goes through expansion and a model trained on partially observed features needs to handle new features in test data without further retraining. The problem is of much significance for dealing with features incrementally collected from different fields. To this end, we propose a new learning paradigm with graph representation and learning. Our framework contains two modules: 1) a backbone network (e.g., feedforward neural nets) as a lower model takes features as input and outputs predicted labels; 2) a graph neural network as an upper model learns to extrapolate embeddings for new features via message passing over a feature-data graph built from observed data. Based on our framework, we design two training strategies, a self-supervised approach and an inductive learning approach, to endow the model with extrapolation ability and alleviate feature-level over-fitting. We also provide theoretical analysis on the generalization error on test data with new features, which dissects the impact of training features and algorithms on generalization performance. Our experiments over several classification datasets and large-scale advertisement click prediction datasets demonstrate that our model can produce effective embeddings for unseen features and significantly outperforms baseline methods that adopt KNN and local aggregation.

[Adversarial Teacher-Student Representation Learning for Domain Generalization](#)

- Fu-En Yang, Yuan-Chia Cheng, Zu-Yun Shiao, Yu-Chiang Frank Wang
- abstract@[open-review](#): Domain generalization (DG) aims to transfer the learning task from a single or multiple source domains to unseen target domains. To extract and leverage the information which exhibits sufficient generalization ability, we propose a simple yet effective approach of Adversarial Teacher-Student Representation Learning, with the goal of deriving the domain generalizable representations via generating and exploring out-of-source data distributions. Our proposed framework advances Teacher-Student learning in an adversarial learning manner, which alternates between knowledge-distillation based representation learning and novel-domain data augmentation. The former progressively updates the teacher network for deriving domain-generalizable representations, while the latter synthesizes data out-of-source yet plausible distributions. Extensive image classification experiments on benchmark datasets in multiple and single source DG settings confirm that, our model exhibits sufficient generalization ability and performs favorably against state-of-the-art DG methods.

[Stochastic bandits with groups of similar arms.](#)

- Fabien Pesquerel, Hassan SABER, Odalric-Ambrym Maillard
- abstract@[open-review](#): We consider a variant of the stochastic multi-armed bandit problem where arms are known to be organized into different groups having the same mean. The groups are unknown but a lower bound $\$q\$$ on their size is known. This situation typically appears when each arm can be described with a list of categorical attributes, and the (unknown) mean reward function only depends on a subset of them, the others being redundant. In this case, $\$q\$$ is linked naturally to the number of attributes considered redundant, and the number of categories of each attribute. For this structured problem of practical relevance, we first derive the asymptotic regret lower bound and corresponding constrained optimization problem. They reveal the achievable regret can be substantially reduced when compared to the unstructured setup, possibly by a factor $\$q\$$. However, solving exactly the exact constrained optimization problem involves a combinatorial problem. We introduce a lower-bound inspired strategy involving a computationally efficient relaxation that is based on a sorting mechanism. We further prove it achieves a lower bound close to the optimal one up to a controlled factor, and achieves an asymptotic regret $\$q\$$ times smaller than the unstructured one. We believe this shows it is a valuable strategy for the practitioner. Last, we illustrate the performance of the considered strategy on numerical experiments involving a large number of arms.

[Tracking Without Re-recognition in Humans and Machines](#)

- Drew Linsley, Girik Malik, Junkyung Kim, Lakshmi Narasimhan Govindarajan, Ennio Mingolla, Thomas Serre
- abstract@[open-review](#): Imagine trying to track one particular fruitfly in a swarm of hundreds. Higher biological visual systems have evolved to track moving objects by relying on both their appearance and their motion trajectories. We investigate if state-of-the-art spatiotemporal deep neural networks are capable of the same. For this, we introduce PathTracker, a synthetic visual challenge that asks human observers and machines to track a target object in the midst of identical-looking "distractor" objects. While humans effortlessly learn PathTracker and generalize to systematic variations in task design, deep networks struggle. To address this limitation, we identify and model circuit mechanisms in biological brains that are implicated in tracking objects based on motion cues. When instantiated as a recurrent network, our circuit model learns to solve PathTracker with a robust visual strategy that rivals human performance and explains a significant proportion of their decision-making on the challenge. We also show that the success of this circuit model extends to object tracking in natural videos. Adding it to a transformer-based architecture for object tracking builds tolerance to visual nuisances that affect object appearance, establishing the new state of the art on the large-scale TrackingNet challenge. Our work highlights the importance of understanding human vision to improve computer vision.

Rethinking conditional GAN training: An approach using geometrically structured latent manifolds

- Sameera Ramasinghe, Moshiur Farazi, Salman H Khan, Nick Barnes, Stephen Gould
- abstract@[open-review](#): Conditional GANs (cGAN), in their rudimentary form, suffer from critical drawbacks such as the lack of diversity in generated outputs and distortion between the latent and output manifolds. Although efforts have been made to improve results, they can suffer from unpleasant side-effects such as the topology mismatch between latent and output spaces. In contrast, we tackle this problem from a geometrical perspective and propose a novel training mechanism that increases both the diversity and the visual quality of a vanilla cGAN, by systematically encouraging a bi-lipschitz mapping between the latent and the output manifolds. We validate the efficacy of our solution on a baseline cGAN (i.e., Pix2Pix) which lacks diversity, and show that by only modifying its training mechanism (i.e., with our proposed Pix2Pix-Geo), one can achieve more diverse and realistic outputs on a broad set of image-to-image translation tasks.

How to transfer algorithmic reasoning knowledge to learn new algorithms?

- Louis-Pascal Xhonneux, Andreea-Ioana Deac, Petar Veličković, Jian Tang
- abstract@[open-review](#): Learning to execute algorithms is a fundamental problem that has been widely studied. Prior work (Veličković et al., 2019) has shown that to enable systematic generalisation on graph algorithms it is critical to have access to the intermediate steps of the program/algorithm. In many reasoning tasks, where algorithmic-style reasoning is important, we only have access to the input and output examples. Thus, inspired by the success of pre-training on similar tasks or data in Natural Language Processing (NLP) and Computer vision, we set out to study how we can transfer algorithmic reasoning knowledge. Specifically, we investigate how we can use algorithms for which we have access to the execution trace to learn to solve similar tasks for which we do not. We investigate two major classes of graph algorithms, parallel algorithms such as breadth-first search and Bellman-Ford and sequential greedy algorithms such as Prims and Dijkstra. Due to the fundamental differences between algorithmic reasoning knowledge and feature extractors such as used in Computer vision or NLP, we hypothesize that standard transfer techniques will not be sufficient to achieve systematic generalisation. To investigate this empirically we create a dataset including 9 algorithms and 3 different graph types. We validate this empirically and show how instead multi-task learning can be used to achieve the transfer of algorithmic reasoning knowledge.

Fast Axiomatic Attribution for Neural Networks

- Robin Hesse, Simone Schaub-Meyer, Stefan Roth
- abstract@[open-review](#): Mitigating the dependence on spurious correlations present in the training dataset is a quickly emerging and important topic of deep learning. Recent approaches include priors on the feature attribution of a deep neural network (DNN) into the training process to reduce the dependence on unwanted features. However, until now one needed to trade off high-quality attributions, satisfying desirable axioms, against the time required to compute them. This in turn either led to long training times or ineffective attribution priors. In this work, we break this trade-off by considering a special class of efficiently axiomatically attributable DNNs for which an axiomatic feature attribution can be computed with only a single forward/backward pass. We formally prove that nonnegatively homogeneous DNNs, here termed \mathcal{X} -DNNs, are efficiently axiomatically attributable and show that they can be effortlessly constructed from a wide range of regular DNNs by simply removing the bias term of each layer. Various experiments demonstrate the advantages of \mathcal{X} -DNNs, beating state-of-the-art generic attribution methods on regular DNNs for training with attribution priors.

OSOA: One-Shot Online Adaptation of Deep Generative Models for Lossless Compression

- Chen Zhang, Shifeng Zhang, Fabio Maria Carlucci, Zhenguo Li
- abstract@[open-review](#): Explicit deep generative models (DGMs), e.g., VAEs and Normalizing Flows, have shown to offer an effective data modelling alternative for lossless compression. However, DGMs themselves normally require large storage space and thus contaminate the advantage brought by accurate data density estimation. To eliminate the requirement of saving separate models for different target datasets, we propose a novel setting that starts from a pretrained deep generative model and compresses the data batches while adapting the model with a dynamical system for only one epoch. We formalise this setting as that of One-Shot Online Adaptation (OSOA) of DGMs for lossless compression and propose a vanilla algorithm under this setting. Experimental results show that vanilla OSOA can save significant time versus training bespoke models and space versus using one model for all targets. With the same adaptation step number or adaptation time, it is shown vanilla OSOA can exhibit better space efficiency, e.g., $\$47\%$ less space, than fine-tuning the pretrained model and saving the fine-tuned model. Moreover, we showcase the potential of OSOA and motivate more sophisticated OSOA algorithms by showing further space or time efficiency with multiple updates per batch and early stopping.

Compressive Visual Representations

- Kuang-Huei Lee, Anurag Arnab, Sergio Guadarrama, John Canny, Ian Fischer
- abstract@[open-review](#): Learning effective visual representations that generalize well without human supervision is a fundamental problem in order to apply Machine Learning to a wide variety of tasks. Recently, two families of self-supervised methods, contrastive learning and latent bootstrapping, exemplified by SimCLR and BYOL respectively, have made significant progress. In this work, we hypothesize that adding explicit information compression to these algorithms yields better and more robust representations. We verify this by developing SimCLR and BYOL formulations compatible with the Conditional Entropy Bottleneck (CEB) objective, allowing us to both measure and control the amount of compression in the learned representation, and observe their impact on downstream tasks. Furthermore, we explore the relationship between Lipschitz continuity and compression, showing a tractable lower bound on the Lipschitz constant of the encoders we learn. As Lipschitz continuity is closely related to robustness, this provides a new explanation for why compressed models are more robust. Our experiments confirm that adding compression to SimCLR and BYOL significantly improves linear evaluation accuracies and model robustness across a wide range of domain shifts. In particular, the compressed version of BYOL achieves 76.0% Top-1 linear evaluation accuracy on ImageNet with ResNet-50, and 78.8% with ResNet-50 2x.

Multi-Armed Bandits with Bounded Arm-Memory: Near-Optimal Guarantees for Best-Arm Identification and Regret Minimization

- Arnab Maiti, Vishakha Patil, Arindam Khan
- abstract@[open-review](#): We study the Stochastic Multi-armed Bandit problem under bounded arm-memory. In this setting, the arms arrive in a stream, and the number of arms that can be stored in the memory at any time, is bounded. The decision-maker can only pull arms that are present in the memory. We address the problem from the perspective of two standard objectives: 1) regret minimization, and 2) best-arm identification. For regret minimization, we settle an important open question by showing an almost tight guarantee. We show $\$O(\Omega(T^{2/3}))$ cumulative regret in expectation for single-pass algorithms for arm-memory size of $\$(n-1)$, where $\$n$ is the number of arms. For best-arm identification, we provide an $\$(\varepsilon, \delta)$ -PAC algorithm with arm memory size of $\$O(\log^* n)$ and $\$O(\frac{n}{\varepsilon^2} \cdot \log(\frac{1}{\delta}))$ optimal sample complexity.

Grounding inductive biases in natural images: invariance stems from variations in data

- Diane Bouchacourt, Mark Ibrahim, Ari Morcos
- abstract@[open-review](#): To perform well on unseen and potentially out-of-distribution samples, it is desirable for machine learning models to have a predictable response with respect to transformations affecting the factors of variation of the input. Here, we study the relative importance of several types

of inductive biases towards such predictable behavior: the choice of data, their augmentations, and model architectures. Invariance is commonly achieved through hand-engineered data augmentation, but do standard data augmentations address transformations that explain variations in real data? While prior work has focused on synthetic data, we attempt here to characterize the factors of variation in a real dataset, ImageNet, and study the invariance of both standard residual networks and the recently proposed vision transformer with respect to changes in these factors. We show standard augmentation relies on a precise combination of translation and scale, with translation recapturing most of the performance improvement--despite the (approximate) translation invariance built in to convolutional architectures, such as residual networks. In fact, we found that scale and translation invariance was similar across residual networks and vision transformer models despite their markedly different architectural inductive biases. We show the training data itself is the main source of invariance, and that data augmentation only further increases the learned invariances. Notably, the invariances learned during training align with the ImageNet factors of variation we found. Finally, we find that the main factors of variation in ImageNet mostly relate to appearance and are specific to each class.

[Directed Graph Contrastive Learning](#)

- Zekun Tong, Yuxuan Liang, Henghui Ding, Yongxing Dai, Xinkle Li, Changhu Wang
- abstract@[open-review](#): Graph Contrastive Learning (GCL) has emerged to learn generalizable representations from contrastive views. However, it is still in its infancy with two concerns: 1) changing the graph structure through data augmentation to generate contrastive views may mislead the message passing scheme, as such graph changing action deprives the intrinsic graph structural information, especially the directional structure in directed graphs; 2) since GCL usually uses predefined contrastive views with hand-picking parameters, it does not take full advantage of the contrastive information provided by data augmentation, resulting in incomplete structure information for models learning. In this paper, we design a directed graph data augmentation method called Laplacian perturbation and theoretically analyze how it provides contrastive information without changing the directed graph structure. Moreover, we present a directed graph contrastive learning framework, which dynamically learns from all possible contrastive views generated by Laplacian perturbation. Then we train it using multi-task curriculum learning to progressively learn from multiple easy-to-difficult contrastive views. We empirically show that our model can retain more structural features of directed graphs than other GCL models because of its ability to provide complete contrastive information. Experiments on various benchmarks reveal our dominance over the state-of-the-art approaches.

[Space-time Mixing Attention for Video Transformer](#)

- Adrian Bulat, Juan Manuel Perez Rua, Swathikiran Sudhakaran, Brais Martinez, Georgios Tzimiropoulos
- abstract@[open-review](#): This paper is on video recognition using Transformers. Very recent attempts in this area have demonstrated promising results in terms of recognition accuracy, yet they have been also shown to induce, in many cases, significant computational overheads due to the additional modelling of the temporal information. In this work, we propose a Video Transformer model the complexity of which scales linearly with the number of frames in the video sequence and hence induces no overhead compared to an image-based Transformer model. To achieve this, our model makes two approximations to the full space-time attention used in Video Transformers: (a) It restricts time attention to a local temporal window and capitalizes on the Transformer's depth to obtain full temporal coverage of the video sequence. (b) It uses efficient space-time mixing to attend jointly spatial and temporal locations without inducing any additional cost on top of a spatial-only attention model. We also show how to integrate 2 very lightweight mechanisms for global temporal-only attention which provide additional accuracy improvements at minimal computational cost. We demonstrate that our model produces very high recognition accuracy on the most popular video recognition datasets while at the same time being significantly more efficient than other Video Transformer models.

[Particle Dual Averaging: Optimization of Mean Field Neural Network with Global Convergence Rate Analysis](#)

- Atsushi Nitanda, Denny Wu, Taiji Suzuki
- abstract@[open-review](#): We propose the particle dual averaging (PDA) method, which generalizes the dual averaging method in convex optimization to the optimization over probability distributions with quantitative runtime guarantee. The algorithm consists of an inner loop and outer loop: the inner loop utilizes the Langevin algorithm to approximately solve for a stationary distribution, which is then optimized in the outer loop. The method can be interpreted as an extension of the Langevin algorithm to naturally handle nonlinear functional on the probability space. An important application of the proposed method is the optimization of neural network in the mean field regime, which is theoretically attractive due to the presence of nonlinear feature learning, but quantitative convergence rate can be challenging to obtain. By adapting finite-dimensional convex optimization theory into the space of measures, we not only establish global convergence of PDA for two-layer mean field neural networks under more general settings and simpler analysis, but also provide quantitative polynomial runtime guarantee. Our theoretical results are supported by numerical simulations on neural networks with reasonable size.

[Learning Tree Interpretation from Object Representation for Deep Reinforcement Learning](#)

- Guiliang Liu, Xiangyu Sun, Oliver Schulte, Pascal Poupart
- abstract@[open-review](#): Interpreting Deep Reinforcement Learning (DRL) models is important to enhance trust and comply with transparency regulations. Existing methods typically explain a DRL model by visualizing the importance of low-level input features with super-pixels, attentions, or saliency maps. Our approach provides an interpretation based on high-level latent object features derived from a disentangled representation. We propose a Represent And Mimic (RAMi) framework for training 1) an identifiable latent representation to capture the independent factors of variation for the objects and 2) a mimic tree that extracts the causal impact of the latent features on DRL action values. To jointly optimize both the fidelity and the simplicity of a mimic tree, we derive a novel Minimum Description Length (MDL) objective based on the Information Bottleneck (IB) principle. Based on this objective, we describe a Monte Carlo Regression Tree Search (MCRTS) algorithm that explores different splits to find the IB-optimal mimic tree. Experiments show that our mimic tree achieves strong approximation performance with significantly fewer nodes than baseline models. We demonstrate the interpretability of our mimic tree by showing latent traversals, decision rules, causal impacts, and human evaluation results.

[Only Train Once: A One-Shot Neural Network Training And Pruning Framework](#)

- Tianyi Chen, Bo Ji, Tianyu Ding, Biyi Fang, Guanyi Wang, Zhihui Zhu, Luming Liang, Yixin Shi, Sheng Yi, Xiao Tu
- abstract@[open-review](#): Structured pruning is a commonly used technique in deploying deep neural networks (DNNs) onto resource-constrained devices. However, the existing pruning methods are usually heuristic, task-specified, and require an extra fine-tuning procedure. To overcome these limitations, we propose a framework that compresses DNNs into slimmer architectures with competitive performances and significant FLOPs reductions by Only-Train-Once (OTO). OTO contains two key steps: (i) we partition the parameters of DNNs into zero-invariant groups, enabling us to prune zero groups without affecting the output; and (ii) to promote zero groups, we then formulate a structured-sparsity optimization problem, and propose a novel optimization algorithm, Half-Space Stochastic Projected Gradient (HSPG), to solve it, which outperforms the standard proximal methods on group sparsity exploration, and maintains comparable convergence. To demonstrate the effectiveness of OTO, we train and compress full models simultaneously from scratch without fine-tuning for inference speedup and parameter reduction, and achieve state-of-the-art results on VGG16 for CIFAR10, ResNet50 for CIFAR10 and Bert for SQuAD and competitive result on ResNet50 for ImageNet. The source code is available at <https://github.com/tianyic/onlytrainonce>.

[Referring Transformer: A One-step Approach to Multi-task Visual Grounding](#)

- Muchen Li, Leonid Sigal

- abstract@[open-review](#): As an important step towards visual reasoning, visual grounding (e.g., phrase localization, referring expression comprehension / segmentation) has been widely explored. Previous approaches to referring expression comprehension (REC) or segmentation (RES) either suffer from limited performance, due to a two-stage setup, or require the designing of complex task-specific one-stage architectures. In this paper, we propose a simple one-stage multi-task framework for visual grounding tasks. Specifically, we leverage a transformer architecture, where two modalities are fused in a visual-lingual encoder. In the decoder, the model learns to generate contextualized lingual queries which are then decoded and used to directly regress the bounding box and produce a segmentation mask for the corresponding referred regions. With this simple but highly contextualized model, we outperform state-of-the-art methods by a large margin on both REC and RES tasks. We also show that a simple pre-training schedule (on an external dataset) further improves the performance. Extensive experiments and ablations illustrate that our model benefits greatly from contextualized information and multi-task training.

[Decoupling the Depth and Scope of Graph Neural Networks](#)

- Hanqing Zeng, Muhan Zhang, Yinglong Xia, Ajitesh Srivastava, Andrey Malevich, Rajgopal Kannan, Viktor Prasanna, Long Jin, Ren Chen
- abstract@[open-review](#): State-of-the-art Graph Neural Networks (GNNs) have limited scalability with respect to the graph and model sizes. On large graphs, increasing the model depth often means exponential expansion of the scope (i.e., receptive field). Beyond just a few layers, two fundamental challenges emerge: 1. degraded expressivity due to oversmoothing, and 2. expensive computation due to neighborhood explosion. We propose a design principle to decouple the depth and scope of GNNs – to generate representation of a target entity (i.e., a node or an edge), we first extract a localized subgraph as the bounded-size scope, and then apply a GNN of arbitrary depth on top of the subgraph. A properly extracted subgraph consists of a small number of critical neighbors, while excluding irrelevant ones. The GNN, no matter how deep it is, smooths the local neighborhood into informative representation rather than oversmoothing the global graph into “white noise”. Theoretically, decoupling improves the GNN expressive power from the perspectives of graph signal processing (GCN), function approximation (GraphSAGE) and topological learning (GIN). Empirically, on seven graphs (with up to 110M nodes) and six backbone GNN architectures, our design achieves significant accuracy improvement with orders of magnitude reduction in computation and hardware cost.

[Fast and Memory Efficient Differentially Private-SGD via JL Projections](#)

- Zhiqi Bu, Sivakanth Gopi, Janardhan Kulkarni, Yin Tat Lee, Hanwen Shen, Uthaipon Tantipongpipat
- abstract@[open-review](#): Differentially Private-SGD (DP-SGD) of Abadi et al. and its variations are the only known algorithms for private training of large scale neural networks. This algorithm requires computation of per-sample gradients norms which is extremely slow and memory intensive in practice. In this paper, we present a new framework to design differentially private optimizers called DP-SGD-JL and DP-Adam-JL. Our approach uses Johnson–Lindenstrauss (JL) projections to quickly approximate the per-sample gradient norms without exactly computing them, thus making the training time and memory requirements of our optimizers closer to that of their non-DP versions. Unlike previous attempts to make DP-SGD faster which work only on a subset of network architectures or use compiler techniques, we propose an algorithmic solution which works for any network in a black-box manner which is the main contribution of this paper. To illustrate this, on IMDb dataset, we train a Recurrent Neural Network (RNN) to achieve good privacy-vs-accuracy tradeoff, while being significantly faster than DP-SGD and with a similar memory footprint as non-private SGD.

[Formalizing Generalization and Adversarial Robustness of Neural Networks to Weight Perturbations](#)

- Yu-Lin Tsai, Chia-Yi Hsu, Chia-Mu Yu, Pin-Yu Chen
- abstract@[open-review](#): Studying the sensitivity of weight perturbation in neural networks and its impacts on model performance, including generalization and robustness, is an active research topic due to its implications on a wide range of machine learning tasks such as model compression, generalization gap assessment, and adversarial attacks. In this paper, we provide the first integral study and analysis for feed-forward neural networks in terms of the robustness in pairwise class margin and its generalization behavior under weight perturbation. We further design a new theory-driven loss function for training generalizable and robust neural networks against weight perturbations. Empirical experiments are conducted to validate our theoretical analysis. Our results offer fundamental insights for characterizing the generalization and robustness of neural networks against weight perturbations.

[Pipeline Combinators for Gradual AutoML](#)

- Guillaume Baudart, Martin Hirzel, Kiran Kate, Parikshit Ram, Avi Shinnar, Jason Tsay
- abstract@[open-review](#): Automated machine learning (AutoML) can make data scientists more productive. But if machine learning is totally automated, that leaves no room for data scientists to apply their intuition. Hence, data scientists often prefer not total but gradual automation, where they control certain choices and AutoML explores the rest. Unfortunately, gradual AutoML is cumbersome with state-of-the-art tools, requiring large non-compositional code changes. More concise compositional code can be achieved with combinators, a powerful concept from functional programming. This paper introduces a small set of orthogonal combinators for composing machine-learning operators into pipelines. It describes a translation scheme from pipelines and associated hyperparameter schemas to search spaces for AutoML optimizers. On that foundation, this paper presents Lale, an open-source sklearn-compatible AutoML library, and evaluates it with a user study.

[Boost Neural Networks by Checkpoints](#)

- Feng Wang, Guoyizhe Wei, Qiao Liu, Jinxiang Ou, xian wei, Hairong Lv
- abstract@[open-review](#): Training multiple deep neural networks (DNNs) and averaging their outputs is a simple way to improve the predictive performance. Nevertheless, the multiplied training cost prevents this ensemble method to be practical and efficient. Several recent works attempt to save and ensemble the checkpoints of DNNs, which only requires the same computational cost as training a single network. However, these methods suffer from either marginal accuracy improvements due to the low diversity of checkpoints or high risk of divergence due to the cyclical learning rates they adopted. In this paper, we propose a novel method to ensemble the checkpoints, where a boosting scheme is utilized to accelerate model convergence and maximize the checkpoint diversity. We theoretically prove that it converges by reducing exponential loss. The empirical evaluation also indicates our proposed ensemble outperforms single model and existing ensembles in terms of accuracy and efficiency. With the same training budget, our method achieves 4.16% lower error on Cifar-100 and 6.96% on Tiny-ImageNet with ResNet-110 architecture. Moreover, the adaptive sample weights in our method make it an effective solution to address the imbalanced class distribution. In the experiments, it yields up to 5.02% higher accuracy over single EfficientNet-B0 on the imbalanced datasets.

[Model Selection for Bayesian Autoencoders](#)

- Ba-Hien Tran, Simone Rossi, Dimitrios Milios, Pietro Michiardi, Edwin V. Bonilla, Maurizio Filippone
- abstract@[open-review](#): We develop a novel method for carrying out model selection for Bayesian autoencoders (BAEs) by means of prior hyper-parameter optimization. Inspired by the common practice of type-II maximum likelihood optimization and its equivalence to Kullback-Leibler divergence minimization, we propose to optimize the distributional sliced-Wasserstein distance (DSWD) between the output of the autoencoder and the empirical data distribution. The advantages of this formulation are that we can estimate the DSWD based on samples and handle high-dimensional problems. We carry out posterior estimation of the BAE parameters via stochastic gradient Hamiltonian Monte Carlo and turn our BAE into a generative model by fitting a flexible Dirichlet mixture model in the latent space. Thanks to this approach, we obtain a powerful alternative to variational autoencoders, which are the preferred choice in modern application of autoencoders for representation learning with uncertainty. We evaluate our approach qualitatively and

quantitatively using a vast experimental campaign on a number of unsupervised learning tasks and show that, in small-data regimes where priors matter, our approach provides state-of-the-art results, outperforming multiple competitive baselines.

[Three Operator Splitting with Subgradients, Stochastic Gradients, and Adaptive Learning Rates](#)

- Alp Yurtsever, Alex Gu, Suvrit Sra
- abstract@[open-review](#): Three Operator Splitting (TOS) (Davis & Yin, 2017) can minimize the sum of multiple convex functions effectively when an efficient gradient oracle or proximal operator is available for each term. This requirement often fails in machine learning applications: (i) instead of full gradients only stochastic gradients may be available; and (ii) instead of proximal operators, using subgradients to handle complex penalty functions may be more efficient and realistic. Motivated by these concerns, we analyze three potentially valuable extensions of TOS. The first two permit using subgradients and stochastic gradients, and are shown to ensure a $\mathcal{O}(1/\sqrt{t})$ convergence rate. The third extension AdapTOS endows TOS with adaptive step-sizes. For the important setting of optimizing a convex loss over the intersection of convex sets AdapTOS attains universal convergence rates, i.e., the rate adapts to the unknown smoothness degree of the objective. We compare our proposed methods with competing methods on various applications.

[Knowledge-Adaptation Priors](#)

- Mohammad Emamiyaz E. Khan, Siddharth Swaroop
- abstract@[open-review](#): Humans and animals have a natural ability to quickly adapt to their surroundings, but machine-learning models, when subjected to changes, often require a complete retraining from scratch. We present Knowledge-adaptation priors (K-priors) to reduce the cost of retraining by enabling quick and accurate adaptation for a wide-variety of tasks and models. This is made possible by a combination of weight and function-space priors to reconstruct the gradients of the past, which recovers and generalizes many existing, but seemingly-unrelated, adaptation strategies. Training with simple first-order gradient methods can often recover the exact retrained model to an arbitrary accuracy by choosing a sufficiently large memory of the past data. Empirical results show that adaptation with K-priors achieves performance similar to full retraining, but only requires training on a handful of past examples.

[Provably efficient multi-task reinforcement learning with model transfer](#)

- Chicheng Zhang, Zhi Wang
- abstract@[open-review](#): We study multi-task reinforcement learning (RL) in tabular episodic Markov decision processes (MDPs). We formulate a heterogeneous multi-player RL problem, in which a group of players concurrently face similar but not necessarily identical MDPs, with a goal of improving their collective performance through inter-player information sharing. We design and analyze a model-based algorithm, and provide gap-dependent and gap-independent regret upper and lower bounds that characterize the intrinsic complexity of the problem.

[Predicting Molecular Conformation via Dynamic Graph Score Matching](#)

- Shitong Luo, Chence Shi, Minkai Xu, Jian Tang
- abstract@[open-review](#): Predicting stable 3D conformations from 2D molecular graphs has been a long-standing challenge in computational chemistry. Recently, machine learning approaches have demonstrated very promising results compared to traditional experimental and physics-based simulation methods. These approaches mainly focus on modeling the local interactions between neighboring atoms on the molecular graphs and overlook the long-range interactions between non-bonded atoms. However, these non-bonded atoms may be proximal to each other in 3D space, and modeling their interactions is of crucial importance to accurately determine molecular conformations, especially for large molecules and multi-molecular complexes. In this paper, we propose a new approach called Dynamic Graph Score Matching (DGSM) for molecular conformation prediction, which models both the local and long-range interactions by dynamically constructing graph structures between atoms according to their spatial proximity during both training and inference. Specifically, the DGSM directly estimates the gradient fields of the logarithm density of atomic coordinates according to the dynamically constructed graphs using score matching methods. The whole framework can be efficiently trained in an end-to-end fashion. Experiments across multiple tasks show that the DGSM outperforms state-of-the-art baselines by a large margin, and it is capable of generating conformations for a broader range of systems such as proteins and multi-molecular complexes.

[When in Doubt: Neural Non-Parametric Uncertainty Quantification for Epidemic Forecasting](#)

- Harshavardhan Kamarthi, Lingkai Kong, Alexander Rodriguez, Chao Zhang, B. Aditya Prakash
- abstract@[open-review](#): Accurate and trustworthy epidemic forecasting is an important problem for public health planning and disease mitigation. Most existing epidemic forecasting models disregard uncertainty quantification, resulting in mis-calibrated predictions. Recent works in deep neural models for uncertainty-aware time-series forecasting also have several limitations; e.g., it is difficult to specify proper priors in Bayesian NNs, while methods like deep ensembling can be computationally expensive. In this paper, we propose to use neural functional processes to fill this gap. We model epidemic time-series with a probabilistic generative process and propose a functional neural process model called EpiFNP, which directly models the probability distribution of the forecast value in a non-parametric way. In EpiFNP, we use a dynamic stochastic correlation graph to model the correlations between sequences, and design different stochastic latent variables to capture functional uncertainty from different perspectives. Our experiments in a real-time flu forecasting setting show that EpiFNP significantly outperforms state-of-the-art models in both accuracy and calibration metrics, up to 2.5x in accuracy and 2.4x in calibration. Additionally, as EpiFNP learns the relations between the current season and similar patterns of historical seasons, it enables interpretable forecasts. Beyond epidemic forecasting, EpiFNP can be of independent interest for advancing uncertainty quantification in deep sequential models for predictive analytics.

[Bounds all around: training energy-based models with bidirectional bounds](#)

- Cong Geng, Jia Wang, Zhiyong Gao, Jes Frellsen, Søren Hauberg
- abstract@[open-review](#): Energy-based models (EBMs) provide an elegant framework for density estimation, but they are notoriously difficult to train. Recent work has established links to generative adversarial networks, where the EBM is trained through a minimax game with a variational value function. We propose a bidirectional bound on the EBM log-likelihood, such that we maximize a lower bound and minimize an upper bound when solving the minimax game. We link one bound to a gradient penalty that stabilizes training, thereby provide grounding for best engineering practice. To evaluate the bounds we develop a new and efficient estimator of the Jacobi-determinant of the EBM generator. We demonstrate that these developments stabilize training and yield high-quality density estimation and sample generation.

[CogView: Mastering Text-to-Image Generation via Transformers](#)

- Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, Jie Tang
- abstract@[open-review](#): Text-to-Image generation in the general domain has long been an open problem, which requires both a powerful generative model and cross-modal understanding. We propose CogView, a 4-billion-parameter Transformer with VQ-VAE tokenizer to advance this problem. We also demonstrate the finetuning strategies for various downstream tasks, e.g. style learning, super-resolution, text-image ranking and fashion design, and

methods to stabilize pretraining, e.g. eliminating NaN losses. CogView achieves the state-of-the-art FID on the blurred MS COCO dataset, outperforming previous GAN-based models and a recent similar work DALL-E.

[Time-independent Generalization Bounds for SGLD in Non-convex Settings](#)

- Tyler Farghly, Patrick Rebeschini
- abstract@[open-review](#): We establish generalization error bounds for stochastic gradient Langevin dynamics (SGLD) with constant learning rate under the assumptions of dissipativity and smoothness, a setting that has received increased attention in the sampling/optimization literature. Unlike existing bounds for SGLD in non-convex settings, ours are time-independent and decay to zero as the sample size increases. Using the framework of uniform stability, we establish time-independent bounds by exploiting the Wasserstein contraction property of the Langevin diffusion, which also allows us to circumvent the need to bound gradients using Lipschitz-like assumptions. Our analysis also supports variants of SGLD that use different discretization methods, incorporate Euclidean projections, or use non-isotropic noise.

[Nonuniform Negative Sampling and Log Odds Correction with Rare Events Data](#)

- HaiYing Wang, Aonan Zhang, Chong Wang
- abstract@[open-review](#): We investigate the issue of parameter estimation with nonuniform negative sampling for imbalanced data. We first prove that, with imbalanced data, the available information about unknown parameters is only tied to the relatively small number of positive instances, which justifies the usage of negative sampling. However, if the negative instances are subsampled to the same level of the positive cases, there is information loss. To maintain more information, we derive the asymptotic distribution of a general inverse probability weighted (IPW) estimator and obtain the optimal sampling probability that minimizes its variance. To further improve the estimation efficiency over the IPW method, we propose a likelihood-based estimator by correcting log odds for the sampled data and prove that the improved estimator has the smallest asymptotic variance among a large class of estimators. It is also more robust to pilot misspecification. We validate our approach on simulated data as well as a real click-through rate dataset with more than 0.3 trillion instances, collected over a period of a month. Both theoretical and empirical results demonstrate the effectiveness of our method.

[Algorithmic stability and generalization of an unsupervised feature selection algorithm](#)

- xinxing wu, Qiang Cheng
- abstract@[open-review](#): Feature selection, as a vital dimension reduction technique, reduces data dimension by identifying an essential subset of input features, which can facilitate interpretable insights into learning and inference processes. Algorithmic stability is a key characteristic of an algorithm regarding its sensitivity to perturbations of input samples. In this paper, we propose an innovative unsupervised feature selection algorithm attaining this stability with provable guarantees. The architecture of our algorithm consists of a feature scorer and a feature selector. The scorer trains a neural network (NN) to globally score all the features, and the selector adopts a dependent sub-NN to locally evaluate the representation abilities for selecting features. Further, we present algorithmic stability analysis and show that our algorithm has a performance guarantee via a generalization error bound. Extensive experimental results on real-world datasets demonstrate superior generalization performance of our proposed algorithm to strong baseline methods. Also, the properties revealed by our theoretical analysis and the stability of our algorithm-selected features are empirically confirmed.

[On learning sparse vectors from mixture of responses](#)

- Nikita Polyanskii
- abstract@[open-review](#): In this paper, we address two learning problems. Suppose a family of $\|\cdot\|$ unknown sparse vectors is fixed, where each vector has at most k non-zero elements. In the first problem, we concentrate on robust learning the supports of all vectors from the family using a sequence of noisy responses. Each response to a query vector shows the sign of the inner product between a randomly chosen vector from the family and the query vector. In the second problem, we aim at designing queries such that all sparse vectors from the family can be approximately reconstructed based on the error-free responses. This learning model was introduced in the work of Gandikota et al., 2020, and these problems can be seen as generalizations of support recovery and approximate recovery problems, well-studied under the framework of 1-bit compressed sensing. As the main contribution of the paper, we prove the existence of learning algorithms for the first problem which work without any assumptions. Under a mild structural assumption on the unknown vectors, we also show the existence of learning algorithms for the second problem and rigorously analyze their query complexity.

[Convergence and Alignment of Gradient Descent with Random Backpropagation Weights](#)

- Ganlin Song, Ruitu Xu, John Lafferty
- abstract@[open-review](#): Stochastic gradient descent with backpropagation is the workhorse of artificial neural networks. It has long been recognized that backpropagation fails to be a biologically plausible algorithm. Fundamentally, it is a non-local procedure---updating one neuron's synaptic weights requires knowledge of synaptic weights or receptive fields of downstream neurons. This limits the use of artificial neural networks as a tool for understanding the biological principles of information processing in the brain. Lillicrap et al. (2016) propose a more biologically plausible "feedback alignment" algorithm that uses random and fixed backpropagation weights, and show promising simulations. In this paper we study the mathematical properties of the feedback alignment procedure by analyzing convergence and alignment for two-layer networks under squared error loss. In the overparameterized setting, we prove that the error converges to zero exponentially fast, and also that regularization is necessary in order for the parameters to become aligned with the random backpropagation weights. Simulations are given that are consistent with this analysis and suggest further generalizations. These results contribute to our understanding of how biologically plausible algorithms might carry out weight learning in a manner different from Hebbian learning, with performance that is comparable with the full non-local backpropagation algorithm.

[Adder Attention for Vision Transformer](#)

- Han Shu, Jiahao Wang, Hanting Chen, Lin Li, Yujiu Yang, Yunhe Wang
- abstract@[open-review](#): Transformer is a new kind of calculation paradigm for deep learning which has shown strong performance on a large variety of computer vision tasks. However, compared with conventional deep models (e.g., convolutional neural networks), vision transformers require more computational resources which cannot be easily deployed on mobile devices. To this end, we present to reduce the energy consumptions using adder neural network (AdderNet). We first theoretically analyze the mechanism of self-attention and the difficulty for applying adder operation into this module. Specifically, the feature diversity, i.e., the rank of attention map using only additions cannot be well preserved. Thus, we develop an adder attention layer that includes an additional identity mapping. With the new operation, vision transformers constructed using additions can also provide powerful feature representations. Experimental results on several benchmarks demonstrate that the proposed approach can achieve highly competitive performance to that of the baselines while achieving an about 2~3× reduction on the energy consumption.

[Reverse engineering learned optimizers reveals known and novel mechanisms](#)

- Niru Maheswaranathan, David Sussillo, Luke Metz, Ruoxi Sun, Jascha Sohl-Dickstein
- abstract@[open-review](#): Learned optimizers are parametric algorithms that can themselves be trained to solve optimization problems. In contrast to baseline optimizers (such as momentum or Adam) that use simple update rules derived from theoretical principles, learned optimizers use flexible, high-dimensional, nonlinear parameterizations. Although this can lead to better performance, their inner workings remain a mystery. How is a given learned

optimizer able to outperform a well tuned baseline? Has it learned a sophisticated combination of existing optimization techniques, or is it implementing completely new behavior? In this work, we address these questions by careful analysis and visualization of learned optimizers. We study learned optimizers trained from scratch on four disparate tasks, and discover that they have learned interpretable behavior, including: momentum, gradient clipping, learning rate schedules, and new forms of learning rate adaptation. Moreover, we show how dynamics and mechanisms inside of learned optimizers orchestrate these computations. Our results help elucidate the previously murky understanding of how learned optimizers work, and establish tools for interpreting future learned optimizers.

[Matching a Desired Causal State via Shift Interventions](#)

- Jiaqi Zhang, Chandler Squires, Caroline Uhler
- abstract@[open-review](#): Transforming a causal system from a given initial state to a desired target state is an important task permeating multiple fields including control theory, biology, and materials science. In causal models, such transformations can be achieved by performing a set of interventions. In this paper, we consider the problem of identifying a shift intervention that matches the desired mean of a system through active learning. We define the Markov equivalence class that is identifiable from shift interventions and propose two active learning strategies that are guaranteed to exactly match a desired mean. We then derive a worst-case lower bound for the number of interventions required and show that these strategies are optimal for certain classes of graphs. In particular, we show that our strategies may require exponentially fewer interventions than the previously considered approaches, which optimize for structure learning in the underlying causal graph. In line with our theoretical results, we also demonstrate experimentally that our proposed active learning strategies require fewer interventions compared to several baselines.

[Unsupervised Noise Adaptive Speech Enhancement by Discriminator-Constrained Optimal Transport](#)

- Hsin-Yi Lin, Huan-Hsin Tseng, Xugang Lu, Yu Tsao
- abstract@[open-review](#): This paper presents a novel discriminator-constrained optimal transport network (DOTN) that performs unsupervised domain adaptation for speech enhancement (SE), which is an essential regression task in speech processing. The DOTN aims to estimate clean references of noisy speech in a target domain, by exploiting the knowledge available from the source domain. The domain shift between training and testing data has been reported to be an obstacle to learning problems in diverse fields. Although rich literature exists on unsupervised domain adaptation for classification, the methods proposed, especially in regressions, remain scarce and often depend on additional information regarding the input data. The proposed DOTN approach tactically fuses the optimal transport (OT) theory from mathematical analysis with generative adversarial frameworks, to help evaluate continuous labels in the target domain. The experimental results on two SE tasks demonstrate that by extending the classical OT formulation, our proposed DOTN outperforms previous adversarial domain adaptation frameworks in a purely unsupervised manner.

[Optimality of variational inference for stochasticblock model with missing links](#)

- Solenne Gaucher, Olga Klopp
- abstract@[open-review](#): Variational methods are extremely popular in the analysis of network data. Statistical guarantees obtained for these methods typically provide asymptotic normality for the problem of estimation of global model parameters under the stochastic block model. In the present work, we consider the case of networks with missing links that is important in application and show that the variational approximation to the maximum likelihood estimator converges at the minimax rate. This provides the first minimax optimal and tractable estimator for the problem of parameter estimation for the stochastic block model with missing links. We complement our results with numerical studies of simulated and real networks, which confirm the advantages of this estimator over current methods.

[Policy Learning Using Weak Supervision](#)

- Jingkang Wang, Hongyi Guo, Zhaowei Zhu, Yang Liu
- abstract@[open-review](#): Most existing policy learning solutions require the learning agents to receive high-quality supervision signals, e.g., rewards in reinforcement learning (RL) or high-quality expert demonstrations in behavioral cloning (BC). These quality supervisions are either infeasible or prohibitively expensive to obtain in practice. We aim for a unified framework that leverages the available cheap weak supervisions to perform policy learning efficiently. To handle this problem, we treat the "weak supervision" as imperfect information coming from a peer agent, and evaluate the learning agent's policy based on a correlated agreement" with the peer agent's policy (instead of simple agreements). Our approach explicitly punishes a policy for overfitting to the weak supervision. In addition to theoretical guarantees, extensive evaluations on tasks including RL with noisy reward, BC with weak demonstrations, and standard policy co-training (RL + BC) show that our method leads to substantial performance improvements, especially when the complexity or the noise of the learning environments is high.

[Chasing Sparsity in Vision Transformers: An End-to-End Exploration](#)

- Tianlong Chen, Yu Cheng, Zhe Gan, Lu Yuan, Lei Zhang, Zhangyang Wang
- abstract@[open-review](#): Vision transformers (ViTs) have recently received explosive popularity, but their enormous model sizes and training costs remain daunting. Conventional post-training pruning often incurs higher training budgets. In contrast, this paper aims to trim down both the training memory overhead and the inference complexity, without sacrificing the achievable accuracy. We carry out the first-of-its-kind comprehensive exploration, on taking a unified approach of integrating sparsity in ViTs "from end to end". Specifically, instead of training full ViTs, we dynamically extract and train sparse subnetworks, while sticking to a fixed small parameter budget. Our approach jointly optimizes model parameters and explores connectivity throughout training, ending up with one sparse network as the final output. The approach is seamlessly extended from unstructured to structured sparsity, the latter by considering to guide the prune-and-grow of self-attention heads inside ViTs. We further co-explore data and architecture sparsity for additional efficiency gains by plugging in a novel learnable token selector to adaptively determine the currently most vital patches. Extensive results on ImageNet with diverse ViT backbones validate the effectiveness of our proposals which obtain significantly reduced computational cost and almost unimpeded generalization. Perhaps most surprisingly, we find that the proposed sparse (co-)training can sometimes \textit{improve} the ViT accuracy rather than compromising it, making sparsity a tantalizing "free lunch". For example, our sparsified DeiT-Small at (\$5\%, \$50\%) sparsity for (data, architecture), improves \$\mathbf{0.28\%}\$ top-1 accuracy, and meanwhile enjoys \$\mathbf{49.32\%}\$ FLOPs and \$\mathbf{4.40\%}\$ running time savings. Our codes are available at <https://github.com/VITA-Group/SViTE>.

[Graphical Models in Heavy-Tailed Markets](#)

- Jose Vinicius de Miranda Cardoso, Jiaxi Ying, Daniel Palomar
- abstract@[open-review](#): Heavy-tailed statistical distributions have long been considered a more realistic statistical model for the data generating process in financial markets in comparison to their Gaussian counterpart. Nonetheless, mathematical nuisances, including nonconvexities, involved in estimating graphs in heavy-tailed settings pose a significant challenge to the practical design of algorithms for graph learning. In this work, we present graph learning estimators based on the Markov random field framework that assume a Student-\$t\$ data generating process. We design scalable numerical algorithms, via the alternating direction method of multipliers, to learn both connected and \$k\$-component graphs along with their theoretical convergence guarantees. The proposed methods outperform state-of-the-art benchmarks in an extensive series of practical experiments with publicly available data from the S&P500 index, foreign exchanges, and cryptocurrencies.

[A Shading-Guided Generative Implicit Model for Shape-Accurate 3D-Aware Image Synthesis](#)

- Xingang Pan, Xudong XU, Chen Change Loy, Christian Theobalt, Bo Dai
- abstract@[open-review](#): The advancement of generative radiance fields has pushed the boundary of 3D-aware image synthesis. Motivated by the observation that a 3D object should look realistic from multiple viewpoints, these methods introduce a multi-view constraint as regularization to learn valid 3D radiance fields from 2D images. Despite the progress, they often fall short of capturing accurate 3D shapes due to the shape-color ambiguity, limiting their applicability in downstream tasks. In this work, we address this ambiguity by proposing a novel shading-guided generative implicit model that is able to learn a starkly improved shape representation. Our key insight is that an accurate 3D shape should also yield a realistic rendering under different lighting conditions. This multi-lighting constraint is realized by modeling illumination explicitly and performing shading with various lighting conditions. Gradients are derived by feeding the synthesized images to a discriminator. To compensate for the additional computational burden of calculating surface normals, we further devise an efficient volume rendering strategy via surface tracking, reducing the training and inference time by 24% and 48%, respectively. Our experiments on multiple datasets show that the proposed approach achieves photorealistic 3D-aware image synthesis while capturing accurate underlying 3D shapes. We demonstrate improved performance of our approach on 3D shape reconstruction against existing methods, and show its applicability on image relighting. Our code is available at <https://github.com/XingangPan/ShadeGAN>.

[XCiT: Cross-Covariance Image Transformers](#)

- Alaaeldin Ali, Hugo Touvron, Mathilde Caron, Piotr Bojanowski, Matthijs Douze, Armand Joulin, Ivan Laptev, Natalia Neverova, Gabriel Synnaeve, Jakob Verbeek, Herve Jegou
- abstract@[open-review](#): Following their success in natural language processing, transformers have recently shown much promise for computer vision. The self-attention operation underlying transformers yields global interactions between all tokens ,i.e. words or image patches, and enables flexible modelling of image data beyond the local interactions of convolutions. This flexibility, however, comes with a quadratic complexity in time and memory, hindering application to long sequences and high-resolution images. We propose a “transposed” version of self-attention that operates across feature channels rather than tokens, where the interactions are based on the cross-covariance matrix between keys and queries. The resulting cross-covariance attention (XCA) has linear complexity in the number of tokens, and allows efficient processing of high-resolution images. Our cross-covariance image transformer (XCiT) is built upon XCA. It combines the accuracy of conventional transformers with the scalability of convolutional architectures. We validate the effectiveness and generality of XCiT by reporting excellent results on multiple vision benchmarks, including image classification and self-supervised feature learning on ImageNet-1k, object detection and instance segmentation on COCO, and semantic segmentation on ADE20k. We will opensource our code and trained models to reproduce the reported results.

[Row-clustering of a Point Process-valued Matrix](#)

- Lihao Yin, Ganggang Xu, Huiyan Sang, Yongtao Guan
- abstract@[open-review](#): Structured point process data harvested from various platforms poses new challenges to the machine learning community. To cluster repeatedly observed marked point processes, we propose a novel mixture model of multi-level marked point processes for identifying potential heterogeneity in the observed data. Specifically, we study a matrix whose entries are marked log-Gaussian Cox processes and cluster rows of such a matrix. An efficient semi-parametric Expectation-Solution (ES) algorithm combined with functional principal component analysis (FPCA) of point processes is proposed for model estimation. The effectiveness of the proposed framework is demonstrated through simulation studies and real data analyses.

[Fine-Grained Neural Network Explanation by Identifying Input Features with Predictive Information](#)

- Yang Zhang, Ashkan Khakzar, Yawei Li, Azade Farshad, Seong Tae Kim, Nassir Navab
- abstract@[open-review](#): One principal approach for illuminating a black-box neural network is feature attribution, i.e. identifying the importance of input features for the network’s prediction. The predictive information of features is recently proposed as a proxy for the measure of their importance. So far, the predictive information is only identified for latent features by placing an information bottleneck within the network. We propose a method to identify features with predictive information in the input domain. The method results in fine-grained identification of input features’ information and is agnostic to network architecture. The core idea of our method is leveraging a bottleneck on the input that only lets input features associated with predictive latent features pass through. We compare our method with several feature attribution methods using mainstream feature attribution evaluation experiments. The code is publicly available.

[Fast Minimum-norm Adversarial Attacks through Adaptive Norm Constraints](#)

- Maura Pintor, Fabio Roli, Wieland Brendel, Battista Biggio
- abstract@[open-review](#): Evaluating adversarial robustness amounts to finding the minimum perturbation needed to have an input sample misclassified. The inherent complexity of the underlying optimization requires current gradient-based attacks to be carefully tuned, initialized, and possibly executed for many computationally-demanding iterations, even if specialized to a given perturbation model. In this work, we overcome these limitations by proposing a fast minimum-norm (FMN) attack that works with different ℓ_p -norm perturbation models ($p=0, 1, 2, \infty$), is robust to hyperparameter choices, does not require adversarial starting points, and converges within few lightweight steps. It works by iteratively finding the sample misclassified with maximum confidence within an ℓ_p -norm constraint of size ϵ , while adapting ϵ to minimize the distance of the current sample to the decision boundary. Extensive experiments show that FMN significantly outperforms existing ℓ_0 , ℓ_1 , and ℓ_∞ -norm attacks in terms of perturbation size, convergence speed and computation time, while reporting comparable performances with state-of-the-art ℓ_2 -norm attacks. Our open-source code is available at: <https://github.com/pralab/Fast-Minimum-Norm-FMN-Attack>.

[Uncertainty Quantification and Deep Ensembles](#)

- Rahul Rahaman, alexandre thiery
- abstract@[open-review](#): Deep Learning methods are known to suffer from calibration issues: they typically produce over-confident estimates. These problems are exacerbated in the low data regime. Although the calibration of probabilistic models is well studied, calibrating extremely over-parametrized models in the low-data regime presents unique challenges. We show that deep-ensembles do not necessarily lead to improved calibration properties. In fact, we show that standard ensembling methods, when used in conjunction with modern techniques such as mixup regularization, can lead to less calibrated models. This text examines the interplay between three of the most simple and commonly used approaches to leverage deep learning when data is scarce: data-augmentation, ensembling, and post-processing calibration methods. We demonstrate that, although standard ensembling techniques certainly help to boost accuracy, the calibration of deep ensembles relies on subtle trade-offs. We also find that calibration methods such as temperature scaling need to be slightly tweaked when used with deep-ensembles and, crucially, need to be executed after the averaging process. Our simulations indicate that, in the low data regime, this simple strategy can halve the Expected Calibration Error (ECE) on a range of benchmark classification problems when compared to standard deep-ensembles.

[Directed Probabilistic Watershed](#)

- Enrique Fita Sanmartin, Sebastian Damrich, Fred A. Hamprecht

- abstract@[open-review](#): The Probabilistic Watershed is a semi-supervised learning algorithm applied on undirected graphs. Given a set of labeled nodes (seeds), it defines a Gibbs probability distribution over all possible spanning forests disconnecting the seeds. It calculates, for every node, the probability of sampling a forest connecting a certain seed with the considered node. We propose the "Directed Probabilistic Watershed", an extension of the Probabilistic Watershed algorithm to directed graphs. Building on the Probabilistic Watershed, we apply the Matrix Tree Theorem for directed graphs and define a Gibbs probability distribution over all incoming directed forests rooted at the seeds. Similar to the undirected case, this turns out to be equivalent to the Directed Random Walker. Furthermore, we show that in the limit case in which the Gibbs distribution has infinitely low temperature, the labeling of the Directed Probabilistic Watershed is equal to the one induced by the incoming directed forest of minimum cost. Finally, for illustration, we compare the empirical performance of the proposed method with other semi-supervised segmentation methods for directed graphs.

[Laplace Redux - Effortless Bayesian Deep Learning](#)

- Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer, Philipp Hennig
- abstract@[open-review](#): Bayesian formulations of deep learning have been shown to have compelling theoretical properties and offer practical functional benefits, such as improved predictive uncertainty quantification and model selection. The Laplace approximation (LA) is a classic, and arguably the simplest family of approximations for the intractable posteriors of deep neural networks. Yet, despite its simplicity, the LA is not as popular as alternatives like variational Bayes or deep ensembles. This may be due to assumptions that the LA is expensive due to the involved Hessian computation, that it is difficult to implement, or that it yields inferior results. In this work we show that these are misconceptions: we (i) review the range of variants of the LA including versions with minimal cost overhead; (ii) introduce "laplace", an easy-to-use software library for PyTorch offering user-friendly access to all major flavors of the LA; and (iii) demonstrate through extensive experiments that the LA is competitive with more popular alternatives in terms of performance, while excelling in terms of computational cost. We hope that this work will serve as a catalyst to a wider adoption of the LA in practical deep learning, including in domains where Bayesian approaches are not typically considered at the moment.

[Hessian Eigenspectra of More Realistic Nonlinear Models](#)

- Zhenyu Liao, Michael W. Mahoney
- abstract@[open-review](#): Given an optimization problem, the Hessian matrix and its eigenspectrum can be used in many ways, ranging from designing more efficient second-order algorithms to performing model analysis and regression diagnostics. When nonlinear models and non-convex problems are considered, strong simplifying assumptions are often made to make Hessian spectral analysis more tractable. This leads to the question of how relevant the conclusions of such analyses are for realistic nonlinear models. In this paper, we exploit tools from random matrix theory to make a precise characterization of the Hessian eigenspectra for a broad family of nonlinear models that extends the classical generalized linear models, without relying on strong simplifying assumptions used previously. We show that, depending on the data properties, the nonlinear response model, and the loss function, the Hessian can have qualitatively different spectral behaviors: of bounded or unbounded support, with single- or multi-bulk, and with isolated eigenvalues on the left- or right-hand side of the main eigenvalue bulk. By focusing on such a simple but nontrivial model, our analysis takes a step forward to unveil the theoretical origin of many visually striking features observed in more realistic machine learning models.

[Explicable Reward Design for Reinforcement Learning Agents](#)

- Rati Devidze, Goran Radanovic, Parameswaran Kamalaruban, Adish Singla
- abstract@[open-review](#): We study the design of explicable reward functions for a reinforcement learning agent while guaranteeing that an optimal policy induced by the function belongs to a set of target policies. By being explicable, we seek to capture two properties: (a) informativeness so that the rewards speed up the agent's convergence, and (b) sparseness as a proxy for ease of interpretability of the rewards. The key challenge is that higher informativeness typically requires dense rewards for many learning tasks, and existing techniques do not allow one to balance these two properties appropriately. In this paper, we investigate the problem from the perspective of discrete optimization and introduce a novel framework, ExpRD, to design explicable reward functions. ExpRD builds upon an informativeness criterion that captures the (sub-)optimality of target policies at different time horizons in terms of actions taken from any given starting state. We provide a mathematical analysis of ExpRD, and show its connections to existing reward design techniques, including potential-based reward shaping. Experimental results on two navigation tasks demonstrate the effectiveness of ExpRD in designing explicable reward functions.

[A Minimalist Approach to Offline Reinforcement Learning](#)

- Scott Fujimoto, Shixiang (Shane) Gu
- abstract@[open-review](#): Offline reinforcement learning (RL) defines the task of learning from a fixed batch of data. Due to errors in value estimation from out-of-distribution actions, most offline RL algorithms take the approach of constraining or regularizing the policy with the actions contained in the dataset. Built on pre-existing RL algorithms, modifications to make an RL algorithm work offline comes at the cost of additional complexity. Offline RL algorithms introduce new hyperparameters and often leverage secondary components such as generative models, while adjusting the underlying RL algorithm. In this paper we aim to make a deep RL algorithm work while making minimal changes. We find that we can match the performance of state-of-the-art offline RL algorithms by simply adding a behavior cloning term to the policy update of an online RL algorithm and normalizing the data. The resulting algorithm is a simple to implement and tune baseline, while more than halving the overall run time by removing the additional computational overheads of previous methods.

[SIMONe: View-Invariant, Temporally-Abstracted Object Representations via Unsupervised Video Decomposition](#)

- Rishabh Kabra, Daniel Zoran, Goker Erdogan, Loic Matthey, Antonia Creswell, Matt Botvinick, Alexander Lerchner, Chris Burgess
- abstract@[open-review](#): To help agents reason about scenes in terms of their building blocks, we wish to extract the compositional structure of any given scene (in particular, the configuration and characteristics of objects comprising the scene). This problem is especially difficult when scene structure needs to be inferred while also estimating the agent's location/viewpoint, as the two variables jointly give rise to the agent's observations. We present an unsupervised variational approach to this problem. Leveraging the shared structure that exists across different scenes, our model learns to infer two sets of latent representations from RGB video input alone: a set of "object" latents, corresponding to the time-invariant, object-level contents of the scene, as well as a set of "frame" latents, corresponding to global time-varying elements such as viewpoint. This factorization of latents allows our model, SIMONe, to represent object attributes in an allocentric manner which does not depend on viewpoint. Moreover, it allows us to disentangle object dynamics and summarize their trajectories as time-abstacted, view-invariant, per-object properties. We demonstrate these capabilities, as well as the model's performance in terms of view synthesis and instance segmentation, across three procedurally generated video datasets.

[Simple Stochastic and Online Gradient Descent Algorithms for Pairwise Learning](#)

- ZHENHUAN YANG, Yunwen Lei, Puyu Wang, Tianbao Yang, Yiming Ying
- abstract@[open-review](#): Pairwise learning refers to learning tasks where the loss function depends on a pair of instances. It instantiates many important machine learning tasks such as bipartite ranking and metric learning. A popular approach to handle streaming data in pairwise learning is an online gradient descent (OGD) algorithm, where one needs to pair the current instance with a buffering set of previous instances with a sufficiently large size and therefore suffers from a scalability issue. In this paper, we propose simple stochastic and online gradient descent methods for pairwise learning. A notable difference from the existing studies is that we only pair the current instance with the previous one in building a gradient direction, which is efficient in both the storage and computational complexity. We develop novel stability results, optimization, and generalization error bounds for both convex and

nonconvex as well as both smooth and nonsmooth problems. We introduce novel techniques to decouple the dependency of models and the previous instance in both the optimization and generalization analysis. Our study resolves an open question on developing meaningful generalization bounds for OGD using a buffering set with a very small fixed size. We also extend our algorithms and stability analysis to develop differentially private SGD algorithms for pairwise learning which significantly improves the existing results.

[User-Level Differentially Private Learning via Correlated Sampling](#)

- Badih Ghazi, Ravi Kumar, Pasin Manurangsi
- abstract@[open-review](#): Most works in learning with differential privacy (DP) have focused on the setting where each user has a single sample. In this work, we consider the setting where each user holds \$m\$ samples and the privacy protection is enforced at the level of each user's data. We show that, in this setting, we may learn with a much fewer number of users. Specifically, we show that, as long as each user receives sufficiently many samples, we can learn any privately learnable class via an (ϵ, δ) -DP algorithm using only $O(\log(1/\delta)/\epsilon)$ users. For ϵ -DP algorithms, we show that we can learn using only $O(\epsilon)(d)$ users even in the local model, where d is the probabilistic representation dimension. In both cases, we show a nearly-matching lower bound on the number of users required. A crucial component of our results is a generalization of global stability [Bun, Livni, Moran, FOCS 2020] that allows the use of public randomness. Under this relaxed notion, we employ a correlated sampling strategy to show that the global stability can be boosted to be arbitrarily close to one, at a polynomial expense in the number of samples.

[Asynchronous Decentralized Online Learning](#)

- Jiyan Jiang, Wenpeng Zhang, Jinjie GU, Wenwu Zhu
- abstract@[open-review](#): Most existing algorithms in decentralized online learning are conducted in the synchronous setting. However, synchronization makes these algorithms suffer from the straggler problem, i.e., fast learners have to wait for slow learners, which significantly reduces such algorithms' overall efficiency. To overcome this problem, we study decentralized online learning in the asynchronous setting, which allows different learners to work at their own pace. We first formulate the framework of Asynchronous Decentralized Online Convex Optimization, which specifies the whole process of asynchronous decentralized online learning using a sophisticated event indexing system. Then we propose the Asynchronous Decentralized Online Gradient-Push (AD-OGP) algorithm, which performs asymmetric gossiping communication and instantaneous model averaging. We further derive a regret bound of AD-OGP, which is a function of the network topology, the levels of processing delays, and the levels of communication delays. Extensive experiments show that AD-OGP runs significantly faster than its synchronous counterpart and also verify the theoretical results.

[Multi-Step Budgeted Bayesian Optimization with Unknown Evaluation Costs](#)

- Raul Astudillo, Daniel Jiang, Maximilian Balandat, Eytan Bakshy, Peter Frazier
- abstract@[open-review](#): Bayesian optimization (BO) is a sample-efficient approach to optimizing costly-to-evaluate black-box functions. Most BO methods ignore how evaluation costs may vary over the optimization domain. However, these costs can be highly heterogeneous and are often unknown in advance in many practical settings, such as hyperparameter tuning of machine learning algorithms or physics-based simulation optimization. Moreover, those few existing methods that acknowledge cost heterogeneity do not naturally accommodate a budget constraint on the total evaluation cost. This combination of unknown costs and a budget constraint introduces a new dimension to the exploration-exploitation trade-off, where learning about the cost incurs a cost itself. Existing methods do not reason about the various trade-offs of this problem in a principled way, leading often to poor performance. We formalize this claim by proving that the expected improvement and the expected improvement per unit of cost, arguably the two most widely used acquisition functions in practice, can be arbitrarily inferior with respect to the optimal non-myopic policy. To overcome the shortcomings of existing approaches, we propose the budgeted multi-step expected improvement, a non-myopic acquisition function that generalizes classical expected improvement to the setting of heterogeneous and unknown evaluation costs. We show that our acquisition function outperforms existing methods in a variety of synthetic and real problems.

[Model-Based Domain Generalization](#)

- Alexander Robey, George J. Pappas, Hamed Hassani
- abstract@[open-review](#): Despite remarkable success in a variety of applications, it is well-known that deep learning can fail catastrophically when presented with out-of-distribution data. Toward addressing this challenge, we consider the \emph{domain generalization} problem, wherein predictors are trained using data drawn from a family of related training domains and then evaluated on a distinct and unseen test domain. We show that under a natural model of data generation and a concomitant invariance condition, the domain generalization problem is equivalent to an infinite-dimensional constrained statistical learning problem; this problem forms the basis of our approach, which we call Model-Based Domain Generalization. Due to the inherent challenges in solving constrained optimization problems in deep learning, we exploit nonconvex duality theory to develop unconstrained relaxations of this statistical problem with tight bounds on the duality gap. Based on this theoretical motivation, we propose a novel domain generalization algorithm with convergence guarantees. In our experiments, we report improvements of up to 30% over state-of-the-art domain generalization baselines on several benchmarks including ColoredMNIST, Camelyon17-WILDS, FMoW-WILDS, and PACS.

[\\$\alpha\\$-IoU: A Family of Power Intersection over Union Losses for Bounding Box Regression](#)

- JIABO HE, Sarah Erfani, Xingjun Ma, James Bailey, Ying Chi, Xian-Sheng Hua
- abstract@[open-review](#): Bounding box (bbox) regression is a fundamental task in computer vision. So far, the most commonly used loss functions for bbox regression are the Intersection over Union (IoU) loss and its variants. In this paper, we generalize existing IoU-based losses to a new family of power IoU losses that have a power IoU term and an additional power regularization term with a single power parameter α . We call this new family of losses the α -IoU losses and analyze properties such as order preservingness and loss/gradient reweighting. Experiments on multiple object detection benchmarks and models demonstrate that α -IoU losses, 1) can surpass existing IoU-based losses by a noticeable performance margin; 2) offer detectors more flexibility in achieving different levels of bbox regression accuracy by modulating α ; and 3) are more robust to small datasets and noisy bboxes.

[Practical Large-Scale Linear Programming using Primal-Dual Hybrid Gradient](#)

- David Applegate, Mateo Diaz, Oliver Hinder, Haihao Lu, Miles Lubin, Brendan O'Donoghue, Warren Schudy
- abstract@[open-review](#): We present PDLP, a practical first-order method for linear programming (LP) that can solve to the high levels of accuracy that are expected in traditional LP applications. In addition, it can scale to very large problems because its core operation is matrix-vector multiplications. PDLP is derived by applying the primal-dual hybrid gradient (PDHG) method, popularized by Chambolle and Pock (2011), to a saddle-point formulation of LP. PDLP enhances PDHG for LP by combining several new techniques with older tricks from the literature; the enhancements include diagonal preconditioning, presolving, adaptive step sizes, and adaptive restarting. PDLP improves the state of the art for first-order methods applied to LP. We compare PDLP with SCS, an ADMM-based solver, on a set of 383 LP instances derived from MIPLIB 2017. With a target of 10^{-8} relative accuracy and 1 hour time limit, PDLP achieves a 6.3x reduction in the geometric mean of solve times and a 4.6x reduction in the number of instances unsolved (from 227 to 49). Furthermore, we highlight standard benchmark instances and a large-scale application (PageRank) where our open-source prototype of PDLP, written in Julia, outperforms a commercial LP solver.

[On the Provable Generalization of Recurrent Neural Networks](#)

- Lifu Wang, Bo Shen, Bo Hu, Xing Cao
- abstract@[open-review](#): Recurrent Neural Network (RNN) is a fundamental structure in deep learning. Recently, some works study the training process of over-parameterized neural networks, and show that over-parameterized networks can learn functions in some notable concept classes with a provable generalization error bound. In this paper, we analyze the training and generalization for RNNs with random initialization, and provide the following improvements over recent works:(1) For a RNN with input sequence $\$x=(X_1,X_2,\dots,X_L)\$$, previous works study to learn functions that are summation of $\$f(\beta^T X_l)\$$ and require normalized conditions that $\|X_l\| \leq \epsilon$ with some very small ϵ depending on the complexity of f . In this paper, using detailed analysis about the neural tangent kernel matrix, we prove a generalization error bound to learn such functions without normalized conditions and show that some notable concept classes are learnable with the numbers of iterations and samples scaling almost-polynomially in the input length L .(2) Moreover, we prove a novel result to learn N -variables functions of input sequence with the form $\$f(\beta^T [X_{\{1\}}, \dots, X_{\{N\}}])\$$, which do not belong to the ``additive'' concept class, i.e., the summation of function $f(X_l)$. And we show that when either N or $\max(l_1, \dots, l_N) - \min(l_1, \dots, l_N)$ is small, $\$f(\beta^T [X_{\{1\}}, \dots, X_{\{N\}}])\$$ will be learnable with the number iterations and samples scaling almost-polynomially in the input length L .

[Differentiable Spline Approximations](#)

- Minsu Cho, Aditya Balu, Ameya Joshi, Anjana Deva Prasad, Biswajit Khara, Soumik Sarkar, Baskar Ganapathysubramanian, Adarsh Krishnamurthy, Chinmay Hegde
- abstract@[open-review](#): The paradigm of differentiable programming has significantly enhanced the scope of machine learning via the judicious use of gradient-based optimization. However, standard differentiable programming methods (such as autodiff) typically require that the machine learning models be differentiable, limiting their applicability. Our goal in this paper is to use a new, principled approach to extend gradient-based optimization to functions well modeled by splines, which encompass a large family of piecewise polynomial models. We derive the form of the (weak) Jacobian of such functions and show that it exhibits a block-sparse structure that can be computed implicitly and efficiently. Overall, we show that leveraging this redesigned Jacobian in the form of a differentiable "layer" in predictive models leads to improved performance in diverse applications such as image segmentation, 3D point cloud reconstruction, and finite element analysis. We also open-source the code at [url{https://github.com/idealab-isu/DSA}](https://github.com/idealab-isu/DSA).

[Rate-Optimal Subspace Estimation on Random Graphs](#)

- Zhixin Zhou, Fan Zhou, Ping Li, Cun-Hui Zhang
- abstract@[open-review](#): We study the theory of random bipartite graph whose adjacency matrix is generated according to a connectivity matrix M . We consider the bipartite graph to be sparse, i.e., the entries of M are upper bounded by certain sparsity parameter. We show that the performance of estimating the connectivity matrix M depends on the sparsity of the graph. We focus on two measurement of performance of estimation: the error of estimating M and the error of estimating the column space of M . In the first case, we consider the operator norm and Frobenius norm of the difference between the estimation and the true connectivity matrix. In the second case, the performance will be measured by the difference between the estimated projection matrix and the true projection matrix in operator norm and Frobenius norm. We will show that the estimators we propose achieve the minimax optimal rate.

[Estimating the Unique Information of Continuous Variables](#)

- Ari Pakman, Amin Nejatbakhsh, Dar Gilboa, Abdullah Makkeh, Luca Mazzucato, Michael Wibral, Elad Schneidman
- abstract@[open-review](#): The integration and transfer of information from multiple sources to multiple targets is a core motive of neural systems. The emerging field of partial information decomposition (PID) provides a novel information-theoretic lens into these mechanisms by identifying synergistic, redundant, and unique contributions to the mutual information between one and several variables. While many works have studied aspects of PID for Gaussian and discrete distributions, the case of general continuous distributions is still uncharted territory. In this work we present a method for estimating the unique information in continuous distributions, for the case of one versus two variables. Our method solves the associated optimization problem over the space of distributions with fixed bivariate marginals by combining copula decompositions and techniques developed to optimize variational autoencoders. We obtain excellent agreement with known analytic results for Gaussians, and illustrate the power of our new approach in several brain-inspired neural models. Our method is capable of recovering the effective connectivity of a chaotic network of rate neurons, and uncovers a complex trade-off between redundancy, synergy and unique information in recurrent networks trained to solve a generalized XOR-task.

[Reliable Causal Discovery with Improved Exact Search and Weaker Assumptions](#)

- Ignavier Ng, Yujia Zheng, Jiji Zhang, Kun Zhang
- abstract@[open-review](#): Many of the causal discovery methods rely on the faithfulness assumption to guarantee asymptotic correctness. However, the assumption can be approximately violated in many ways, leading to sub-optimal solutions. Although there is a line of research in Bayesian network structure learning that focuses on weakening the assumption, such as exact search methods with well-defined score functions, they do not scale well to large graphs. In this work, we introduce several strategies to improve the scalability of exact score-based methods in the linear Gaussian setting. In particular, we develop a super-structure estimation method based on the support of inverse covariance matrix which requires assumptions that are strictly weaker than faithfulness, and apply it to restrict the search space of exact search. We also propose a local search strategy that performs exact search on the local clusters formed by each variable and its neighbors within two hops in the super-structure. Numerical experiments validate the efficacy of the proposed procedure, and demonstrate that it scales up to hundreds of nodes with a high accuracy.

[Node Dependent Local Smoothing for Scalable Graph Learning](#)

- Wentao Zhang, Mingyu Yang, Zeang Sheng, Yang Li, Wen Ouyang, Yangyu Tao, Zhi Yang, Bin CUI
- abstract@[open-review](#): Recent works reveal that feature or label smoothing lies at the core of Graph Neural Networks (GNNs). Concretely, they show feature smoothing combined with simple linear regression achieves comparable performance with the carefully designed GNNs, and a simple MLP model with label smoothing of its prediction can outperform the vanilla GCN. Though an interesting finding, smoothing has not been well understood, especially regarding how to control the extent of smoothness. Intuitively, too small or too large smoothing iterations may cause under-smoothing or over-smoothing and can lead to sub-optimal performance. Moreover, the extent of smoothness is node-specific, depending on its degree and local structure. To this end, we propose a novel algorithm called node-dependent local smoothing (NDLS), which aims to control the smoothness of every node by setting a node-specific smoothing iteration. Specifically, NDLS computes influence scores based on the adjacency matrix and selects the iteration number by setting a threshold on the scores. Once selected, the iteration number can be applied to both feature smoothing and label smoothing. Experimental results demonstrate that NDLS enjoys high accuracy -- state-of-the-art performance on node classifications tasks, flexibility -- can be incorporated with any models, scalability and efficiency -- can support large scale graphs with fast training.

[Parallel and Efficient Hierarchical k-Median Clustering](#)

- Vincent Cohen-Addad, Silvio Lattanzi, Ashkan Norouzi-Fard, Christian Sohler, Ola Svensson

- abstract@[open-review](#): As a fundamental unsupervised learning task, hierarchical clustering has been extensively studied in the past decade. In particular, standard metric formulations as hierarchical k -center, k -means, and k -median received a lot of attention and the problems have been studied extensively in different models of computation. Despite all this interest, not many efficient parallel algorithms are known for these problems. In this paper we introduce a new parallel algorithm for the Euclidean hierarchical k -median problem that, when using machines with memory s (for $s \in \Omega(\log^2(n + \Delta + d))$), outputs a hierarchical clustering such that for every fixed value of k the cost of the solution is at most an $O(\min\{d, \log n \} \log \Delta)$ factor larger in expectation than that of an optimal solution. Furthermore, we also get that for all k simultaneously the cost of the solution is at most an $O(\min\{d, \log n \} \log \Delta \log (\Delta d n))$ factor bigger than the corresponding optimal solution. The algorithm requires in $O(\log_s(nd \log(n + \Delta)))$ rounds. Here d is the dimension of the data set and Δ is the ratio between the maximum and minimum distance of two points in the input dataset. To the best of our knowledge, this is the first *parallel* algorithm for the hierarchical k -median problem with theoretical guarantees. We further complement our theoretical results with an empirical study of our algorithm that shows its effectiveness in practice.

[Human-Adversarial Visual Question Answering](#)

- Sasha Sheng, Amanpreet Singh, Vedanuj Goswami, Jose Magana, Tristan Thrush, Wojciech Galuba, Devi Parikh, Douwe Kiela
- abstract@[open-review](#): Performance on the most commonly used Visual Question Answering dataset (VQA v2) is starting to approach human accuracy. However, in interacting with state-of-the-art VQA models, it is clear that the problem is far from being solved. In order to stress test VQA models, we benchmark them against human-adversarial examples. Human subjects interact with a state-of-the-art VQA model, and for each image in the dataset, attempt to find a question where the model's predicted answer is incorrect. We find that a wide range of state-of-the-art models perform poorly when evaluated on these examples. We conduct an extensive analysis of the collected adversarial examples and provide guidance on future research directions. We hope that this Adversarial VQA (AdVQA) benchmark can help drive progress in the field and advance the state of the art.

[Across-animal odor decoding by probabilistic manifold alignment](#)

- Pedro Herrero-Vidal, Dmitry Rinberg, Cristina Savin
- abstract@[open-review](#): Identifying the common structure of neural dynamics across subjects is key for extracting unifying principles of brain computation and for many brain machine interface applications. Here, we propose a novel probabilistic approach for aligning stimulus-evoked responses from multiple animals in a common low dimensional manifold and use hierarchical inference to identify which stimulus drives neural activity in any given trial. Our probabilistic decoder is robust to a range of features of the neural responses and significantly outperforms existing neural alignment procedures. When applied to recordings from the mouse olfactory bulb, our approach reveals low-dimensional population dynamics that are odor specific and have consistent structure across animals. Thus, our decoder can be used for increasing the robustness and scalability of neural-based chemical detection.

[Excess Capacity and Backdoor Poisoning](#)

- Naren Manoj, Avrim Blum
- abstract@[open-review](#): A backdoor data poisoning attack is an adversarial attack wherein the attacker injects several watermarked, mislabeled training examples into a training set. The watermark does not impact the test-time performance of the model on typical data; however, the model reliably errs on watermarked examples. To gain a better foundational understanding of backdoor data poisoning attacks, we present a formal theoretical framework within which one can discuss backdoor data poisoning attacks for classification problems. We then use this to analyze important statistical and computational issues surrounding these attacks. On the statistical front, we identify a parameter we call the memorization capacity that captures the intrinsic vulnerability of a learning problem to a backdoor attack. This allows us to argue about the robustness of several natural learning problems to backdoor attacks. Our results favoring the attacker involve presenting explicit constructions of backdoor attacks, and our robustness results show that some natural problem settings cannot yield successful backdoor attacks. From a computational standpoint, we show that under certain assumptions, adversarial training can detect the presence of backdoors in a training set. We then show that under similar assumptions, two closely related problems we call backdoor filtering and robust generalization are nearly equivalent. This implies that it is both asymptotically necessary and sufficient to design algorithms that can identify watermarked examples in the training set in order to obtain a learning algorithm that both generalizes well to unseen data and is robust to backdoors.

[A Convergence Analysis of Gradient Descent on Graph Neural Networks](#)

- Pranjal Awasthi, Abhimanyu Das, Sreenivas Gollapudi
- abstract@[open-review](#): Graph Neural Networks (GNNs) are a powerful class of architectures for solving learning problems on graphs. While many variants of GNNs have been proposed in the literature and have achieved strong empirical performance, their theoretical properties are less well understood. In this work we study the convergence properties of the gradient descent algorithm when used to train GNNs. In particular, we consider the realizable setting where the data is generated from a network with unknown weights and our goal is to study conditions under which gradient descent on a GNN architecture can recover near optimal solutions. While such analysis has been performed in recent years for other architectures such as fully connected feed-forward networks, the message passing nature of the updates in a GNN poses a new challenge in understanding the nature of the gradient descent updates. We take a step towards overcoming this by proving that for the case of deep linear GNNs gradient descent provably recovers solutions up to error ϵ in $O(\text{log}(1/\epsilon))$ iterations, under natural assumptions on the data distribution. Furthermore, for the case of one-round GNNs with ReLU activations, we show that gradient descent provably recovers solutions up to error ϵ in $O(\frac{1}{\epsilon} \cdot \frac{\log(1/\epsilon)}{\epsilon})$ iterations.

[Differentiable rendering with perturbed optimizers](#)

- Quentin Le Lidec, Ivan Laptev, Cordelia Schmid, Justin Carpentier
- abstract@[open-review](#): Reasoning about 3D scenes from their 2D image projections is one of the core problems in computer vision. Solutions to this inverse and ill-posed problem typically involve a search for models that best explain observed image data. Notably, images depend both on the properties of observed scenes and on the process of image formation. Hence, if optimization techniques should be used to explain images, it is crucial to design differentiable functions for the projection of 3D scenes into images, also known as differentiable rendering. Previous approaches to differentiable rendering typically replace non-differentiable operations by smooth approximations, impacting the subsequent 3D estimation. In this paper, we take a more general approach and study differentiable renderers through the prism of randomized optimization and the related notion of perturbed optimizers. In particular, our work highlights the link between some well-known differentiable renderer formulations and randomly smoothed optimizers, and introduces differentiable perturbed renderers. We also propose a variance reduction mechanism to alleviate the computational burden inherent to perturbed optimizers and introduce an adaptive scheme to automatically adjust the smoothing parameters of the rendering process. We apply our method to 3D scene reconstruction and demonstrate its advantages on the tasks of 6D pose estimation and 3D mesh reconstruction. By providing informative gradients that can be used as a strong supervisory signal, we demonstrate the benefits of perturbed renderers to obtain more accurate solutions when compared to the state-of-the-art alternatives using smooth gradient approximations.

[BCORLE\(\\$\lambda\\$\): An Offline Reinforcement Learning and Evaluation Framework for Coupons Allocation in E-commerce Market](#)

- Yang Zhang, Bo Tang, Qingyu Yang, Dou An, Hongyin Tang, Chenyang Xi, Xueying LI, Feiyu Xiong

- abstract@[open-review](#): Coupons allocation is an important tool for enterprises to increase the activity and loyalty of users on the e-commerce market. One fundamental problem related is how to allocate coupons within a fixed budget while maximizing users' retention on the e-commerce platform. The online e-commerce environment is complicated and ever changing, so it requires the coupons allocation policy learning can quickly adapt to the changes of the company's business strategy. Unfortunately, existing studies with a huge computation overhead can hardly satisfy the requirements of real-time and fast-response in the real world. Specifically, the problem of coupons allocation within a fixed budget is usually formulated as a Lagrangian problem. Existing solutions need to re-learn the policy once the value of Lagrangian multiplier variable λ is updated, causing a great computation overhead. Besides, a mature e-commerce market often faces tens of millions of users and dozens of types of coupons which construct the huge policy space, further increasing the difficulty of solving the problem. To tackle with above problems, we propose a budget constrained offline reinforcement learning and evaluation with λ -generalization (BCORLE(λ)) framework. The proposed method can help enterprises develop a coupons allocation policy which greatly improves users' retention rate on the platform while ensuring the cost does not exceed the budget. Specifically, λ -generalization method is proposed to lead the policy learning process can be executed according to different λ values adaptively, avoiding re-learning new policies from scratch. Thus the computation overhead is greatly reduced. Further, a novel offline reinforcement learning method and an off-policy evaluation algorithm are proposed for policy learning and policy evaluation, respectively. Finally, experiments on the simulation platform and real-world e-commerce market validate the effectiveness of our approach.

Nested Variational Inference

- Heiko Zimmermann, Hao Wu, Babak Esmaeili, Jan-Willem van de Meent
- abstract@[open-review](#): We develop nested variational inference (NVI), a family of methods that learn proposals for nested importance samplers by minimizing an forward or reverse KL divergence at each level of nesting. NVI is applicable to many commonly-used importance sampling strategies and provides a mechanism for learning intermediate densities, which can serve as heuristics to guide the sampler. Our experiments apply NVI to (a) sample from a multimodal distribution using a learned annealing path (b) learn heuristics that approximate the likelihood of future observations in a hidden Markov model and (c) to perform amortized inference in hierarchical deep generative models. We observe that optimizing nested objectives leads to improved sample quality in terms of log average weight and effective sample size.

Exponential Bellman Equation and Improved Regret Bounds for Risk-Sensitive Reinforcement Learning

- Yingjie Fei, Zhuoran Yang, Yudong Chen, Zhaoran Wang
- abstract@[open-review](#): We study risk-sensitive reinforcement learning (RL) based on the entropic risk measure. Although existing works have established non-asymptotic regret guarantees for this problem, they leave open an exponential gap between the upper and lower bounds. We identify the deficiencies in existing algorithms and their analysis that result in such a gap. To remedy these deficiencies, we investigate a simple transformation of the risk-sensitive Bellman equations, which we call the exponential Bellman equation. The exponential Bellman equation inspires us to develop a novel analysis of Bellman backup procedures in risk-sensitive RL algorithms, and further motivates the design of a novel exploration mechanism. We show that these analytic and algorithmic innovations together lead to improved regret upper bounds over existing ones.

On sensitivity of meta-learning to support data

- Mayank Agarwal, Mikhail Yurochkin, Yuekai Sun
- abstract@[open-review](#): Meta-learning algorithms are widely used for few-shot learning. For example, image recognition systems that readily adapt to unseen classes after seeing only a few labeled examples. Despite their success, we show that modern meta-learning algorithms are extremely sensitive to the data used for adaptation, i.e. support data. In particular, we demonstrate the existence of (unaltered, in-distribution, natural) images that, when used for adaptation, yield accuracy as low as 4% or as high as 95% on standard few-shot image classification benchmarks. We explain our empirical findings in terms of class margins, which in turn suggests that robust and safe meta-learning requires larger margins than supervised learning.

On Large-Cohort Training for Federated Learning

- Zachary Charles, Zachary Garrett, Zhouyuan Huo, Sergei Shmulyan, Virginia Smith
- abstract@[open-review](#): Federated learning methods typically learn a model by iteratively sampling updates from a population of clients. In this work, we explore how the number of clients sampled at each round (the cohort size) impacts the quality of the learned model and the training dynamics of federated learning algorithms. Our work poses three fundamental questions. First, what challenges arise when trying to scale federated learning to larger cohorts? Second, what parallels exist between cohort sizes in federated learning, and batch sizes in centralized learning? Last, how can we design federated learning methods that effectively utilize larger cohort sizes? We give partial answers to these questions based on extensive empirical evaluation. Our work highlights a number of challenges stemming from the use of larger cohorts. While some of these (such as generalization issues and diminishing returns) are analogs of large-batch training challenges, others (including catastrophic training failures and fairness concerns) are unique to federated learning.

Generic Neural Architecture Search via Regression

- Yuhong Li, Cong Hao, Pan Li, Jinjun Xiong, Deming Chen
- abstract@[open-review](#): Most existing neural architecture search (NAS) algorithms are dedicated to and evaluated by the downstream tasks, e.g., image classification in computer vision. However, extensive experiments have shown that, prominent neural architectures, such as ResNet in computer vision and LSTM in natural language processing, are generally good at extracting patterns from the input data and perform well on different downstream tasks. In this paper, we attempt to answer two fundamental questions related to NAS. (1) Is it necessary to use the performance of specific downstream tasks to evaluate and search for good neural architectures? (2) Can we perform NAS effectively and efficiently while being agnostic to the downstream tasks? To answer these questions, we propose a novel and generic NAS framework, termed Generic NAS (GenNAS). GenNAS does not use task-specific labels but instead adopts regression on a set of manually designed synthetic signal bases for architecture evaluation. Such a self-supervised regression task can effectively evaluate the intrinsic power of an architecture to capture and transform the input signal patterns, and allow more sufficient usage of training samples. Extensive experiments across 13 CNN search spaces and one NLP space demonstrate the remarkable efficiency of GenNAS using regression, in terms of both evaluating the neural architectures (quantified by the ranking correlation Spearman's rho between the approximated performances and the downstream task performances) and the convergence speed for training (within a few seconds). For example, on NAS-Bench-101, GenNAS achieves 0.85 rho while the existing efficient methods only achieve 0.38. We then propose an automatic task search to optimize the combination of synthetic signals using limited downstream-task-specific labels, further improving the performance of GenNAS. We also thoroughly evaluate GenNAS's generality and end-to-end NAS performance on all search spaces, which outperforms almost all existing works with significant speedup. For example, on NASBench-201, GenNAS can find near-optimal architectures within 0.3 GPU hour.

The best of both worlds: stochastic and adversarial episodic MDPs with unknown transition

- Tiancheng Jin, Longbo Huang, Haipeng Luo
- abstract@[open-review](#): We consider the best-of-both-worlds problem for learning an episodic Markov Decision Process through T episodes, with the goal of achieving $\widetilde{O}(\sqrt{T})$ regret when the losses are adversarial and simultaneously $O(\log T)$ regret when the losses are (almost) stochastic. Recent work by [Jin and Luo, 2020] achieves this goal when the fixed transition is known, and leaves the case of unknown transition as a major open question. In this work, we resolve this open problem by using the same Follow-the-Regularized-Leader (FTRL) framework together with a set of new techniques. Specifically, we first propose a loss-shifting trick in the FTRL analysis, which greatly simplifies the approach of [Jin

and Luo, 2020] and already improves their results for the known transition case. Then, we extend this idea to the unknown transition case and develop a novel analysis which upper bounds the transition estimation error by the regret itself in the stochastic setting, a key property to ensure $\mathcal{O}(\log T)$ regret.

[Private learning implies quantum stability](#)

- Yihui Quek, Srinivasan Arunachalam, John A Smolin
- abstract@[open-review](#): Learning an unknown n-qubit quantum state rho is a fundamental challenge in quantum computing. Information-theoretically, it is known that tomography requires exponential in n many copies of rho to estimate its entries. Motivated by learning theory, Aaronson et al. introduced many (weaker) learning models: the PAC model of learning states (Proceedings of Royal Society A'07), shadow tomography (STOC'18) for learning "shadows" of a state, a model that also requires learners to be differentially private (STOC'19) and the online model of learning states (NeurIPS'18). In these models it was shown that an unknown state can be learned approximately" using linear in n many copies of rho. But is there any relationship between these models? In this paper we prove a sequence of (information-theoretic) implications from differentially-private PAC learning to online learning and then to quantum stability. Our main result generalizes the recent work of Bun, Livni and Moran (Journal of the ACM'21) who showed that finite Littlestone dimension (of Boolean-valued concept classes) implies PAC learnability in the (approximate) differentially private (DP) setting. We first consider their work in the real-valued setting and further extend to their techniques to the setting of learning quantum states. Key to our results is our generic quantum online learner, Robust Standard Optimal Algorithm (RSOA), which is robust to adversarial imprecision. We then show information-theoretic implications between DP learning quantum states in the PAC model, learnability of quantum states in the one-way communication model, online learning of quantum states, quantum stability (which is our conceptual contribution), various combinatorial parameters and give further applications to gentle shadow tomography and noisy quantum state learning.

[Interesting Object, Curious Agent: Learning Task-Agnostic Exploration](#)

- Simone Parisi, Victoria Dean, Deepak Pathak, Abhinav Gupta
- abstract@[open-review](#): Common approaches for task-agnostic exploration learn tabula-rasa --the agent assumes isolated environments and no prior knowledge or experience. However, in the real world, agents learn in many environments and always come with prior experiences as they explore new ones. Exploration is a lifelong process. In this paper, we propose a paradigm change in the formulation and evaluation of task-agnostic exploration. In this setup, the agent first learns to explore across many environments without any extrinsic goal in a task-agnostic manner. Later on, the agent effectively transfers the learned exploration policy to better explore new environments when solving tasks. In this context, we evaluate several baseline exploration strategies and present a simple yet effective approach to learning task-agnostic exploration policies. Our key idea is that there are two components of exploration: (1) an agent-centric component encouraging exploration of unseen parts of the environment based on an agent's belief; (2) an environment-centric component encouraging exploration of inherently interesting objects. We show that our formulation is effective and provides the most consistent exploration across several training-testing environment pairs. We also introduce benchmarks and metrics for evaluating task-agnostic exploration strategies. The source code is available at <https://github.com/sparisi/cbet/>.

[SimiGrad: Fine-Grained Adaptive Batching for Large Scale Training using Gradient Similarity Measurement](#)

- Heyang Qin, Samyam Rajbhandari, Olatunji Ruwase, Feng Yan, Lei Yang, Yuxiong He
- abstract@[open-review](#): Large scale training requires massive parallelism to finish the training within a reasonable amount of time. To support massive parallelism, large batch training is the key enabler but often at the cost of generalization performance. Existing works explore adaptive batching or hand-tuned static large batching, in order to strike a balance between the computational efficiency and the performance. However, these methods can provide only coarse-grained adaption (e.g., at a epoch level) due to the intrinsic expensive calculation or hand tuning requirements. In this paper, we propose a fully automated and lightweight adaptive batching methodology to enable fine-grained batch size adaption (e.g., at a mini-batch level) that can achieve state-of-the-art performance with record breaking batch sizes. The core component of our method is a lightweight yet efficient representation of the critical gradient noise information. We open-source the proposed methodology by providing a plugin tool that supports mainstream machine learning frameworks. Extensive evaluations on popular benchmarks (e.g., CIFAR10, ImageNet, and BERT-Large) demonstrate that the proposed methodology outperforms state-of-the-art methodologies using adaptive batching approaches or hand-tuned static strategies in both performance and batch size. Particularly, we achieve a new state-of-the-art batch size of 78k in BERT-Large pretraining with SQuAD score 90.69 compared to 90.58 reported in previous state-of-the-art with 59k batch size.

[Variational Inference for Continuous-Time Switching Dynamical Systems](#)

- Lukas Köhs, Bastian Alt, Heinz Koepll
- abstract@[open-review](#): Switching dynamical systems provide a powerful, interpretable modeling framework for inference in time-series data in, e.g., the natural sciences or engineering applications. Since many areas, such as biology or discrete-event systems, are naturally described in continuous time, we present a model based on a Markov jump process modulating a subordinated diffusion process. We provide the exact evolution equations for the prior and posterior marginal densities, the direct solutions of which are however computationally intractable. Therefore, we develop a new continuous-time variational inference algorithm, combining a Gaussian process approximation on the diffusion level with posterior inference for Markov jump processes. By minimizing the path-wise Kullback-Leibler divergence we obtain (i) Bayesian latent state estimates for arbitrary points on the real axis and (ii) point estimates of unknown system parameters, utilizing variational expectation maximization. We extensively evaluate our algorithm under the model assumption and for real-world examples.

[Implicit Regularization in Matrix Sensing via Mirror Descent](#)

- Fan Wu, Patrick Rebeschini
- abstract@[open-review](#): We study discrete-time mirror descent applied to the unregularized empirical risk in matrix sensing. In both the general case of rectangular matrices and the particular case of positive semidefinite matrices, a simple potential-based analysis in terms of the Bregman divergence allows us to establish convergence of mirror descent--with different choices of the mirror maps--to a matrix that, among all global minimizers of the empirical risk, minimizes a quantity explicitly related to the nuclear norm, the Frobenius norm, and the von Neumann entropy. In both cases, this characterization implies that mirror descent, a first-order algorithm minimizing the unregularized empirical risk, recovers low-rank matrices under the same set of assumptions that are sufficient to guarantee recovery for nuclear-norm minimization. When the sensing matrices are symmetric and commute, we show that gradient descent with full-rank factorized parametrization is a first-order approximation to mirror descent, in which case we obtain an explicit characterization of the implicit bias of gradient flow as a by-product.

[STORM+: Fully Adaptive SGD with Recursive Momentum for Nonconvex Optimization](#)

- Kfir Levy, Ali Kavis, Volkan Cevher
- abstract@[open-review](#): In this work we investigate stochastic non-convex optimization problems where the objective is an expectation over smooth loss functions, and the goal is to find an approximate stationary point. The most popular approach to handling such problems is variance reduction techniques, which are also known to obtain tight convergence rates, matching the lower bounds in this case. Nevertheless, these techniques require a careful maintenance of anchor points in conjunction with appropriately selected ``mega-batchesizes''. This leads to a challenging hyperparameter tuning problem, that weakens their practicality. Recently, [Cutkosky and Orabona, 2019] have shown that one can employ recursive momentum in order to avoid the use of

anchor points and large batchsizes, and still obtain the optimal rate for this setting. Yet, their method called \$\\rm{STORM}\$ crucially relies on the knowledge of the smoothness, as well a bound on the gradient norms. In this work we propose \$\\rm{STORM}^+\$, a new method that is completely parameter-free, does not require large batch-sizes, and obtains the optimal \$O(1/T^{1/3})\$ rate for finding an approximate stationary point. Our work builds on the \$\\rm{STORM}\$ algorithm, in conjunction with a novel approach to adaptively set the learning rate and momentum parameters.

[Skipping the Frame-Level: Event-Based Piano Transcription With Neural Semi-CRFs](#)

- Yujia Yan, Frank Cwitkowitz, Zhiyao Duan
- abstract@[open-review](#): Piano transcription systems are typically optimized to estimate pitch activity at each frame of audio. They are often followed by carefully designed heuristics and post-processing algorithms to estimate note events from the frame-level predictions. Recent methods have also framed piano transcription as a multi-task learning problem, where the activation of different stages of a note event are estimated independently. These practices are not well aligned with the desired outcome of the task, which is the specification of note intervals as holistic events, rather than the aggregation of disjoint observations. In this work, we propose a novel formulation of piano transcription, which is optimized to directly predict note events. Our method is based on Semi-Markov Conditional Random Fields (semi-CRF), which produce scores for intervals rather than individual frames. When formulating piano transcription in this way, we eliminate the need to rely on disjoint frame-level estimates for different stages of a note event. We conduct experiments on the MAESTRO dataset and demonstrate that the proposed model surpasses the current state-of-the-art for piano transcription. Our results suggest that the semi-CRF output layer, while still quadratic in complexity, is a simple, fast and well-performing solution for event-based prediction, and may lead to similar success in other areas which currently rely on frame-level estimates.

[Deep Learning on a Data Diet: Finding Important Examples Early in Training](#)

- Mansheej Paul, Surya Ganguli, Gintare Karolina Dziugaite
- abstract@[open-review](#): Recent success in deep learning has partially been driven by training increasingly overparametrized networks on ever larger datasets. It is therefore natural to ask: how much of the data is superfluous, which examples are important for generalization, and how do we find them? In this work, we make the striking observation that, in standard vision datasets, simple scores averaged over several weight initializations can be used to identify important examples very early in training. We propose two such scores—the Gradient Normed (GraNd) and the Error L2-Norm (EL2N) scores—and demonstrate their efficacy on a range of architectures and datasets by pruning significant fractions of training data without sacrificing test accuracy. In fact, using EL2N scores calculated a few epochs into training, we can prune half of the CIFAR10 training set while slightly improving test accuracy. Furthermore, for a given dataset, EL2N scores from one architecture or hyperparameter configuration generalize to other configurations. Compared to recent work that prunes data by discarding examples that are rarely forgotten over the course of training, our scores use only local information early in training. We also use our scores to detect noisy examples and study training dynamics through the lens of important examples—we investigate how the data distribution shapes the loss surface and identify subspaces of the model’s data representation that are relatively stable over training.

[BNS: Building Network Structures Dynamically for Continual Learning](#)

- Qi Qin, Wenpeng Hu, Han Peng, Dongyan Zhao, Bing Liu
- abstract@[open-review](#): Continual learning (CL) of a sequence of tasks is often accompanied with the catastrophic forgetting(CF) problem. Existing research has achieved remarkable results in overcoming CF, especially for task continual learning. However, limited work has been done to achieve another important goal of CL,knowledge transfer.In this paper, we propose a technique (called BNS) to do both. The novelty of BNS is that it dynamically builds a network to learn each new task to overcome CF and to transfer knowledge across tasks at the same time. Experimental results show that when the tasks are different (with little shared knowledge), BNS can already outperform the state-of-the-art baselines. When the tasks are similar and have shared knowledge, BNS outperforms the baselines substantially by a large margin due to its knowledge transfer capability.

[Auditing Black-Box Prediction Models for Data Minimization Compliance](#)

- Bashir Rastegarpanah, Krishna Gummadi, Mark Crovella
- abstract@[open-review](#): In this paper, we focus on auditing black-box prediction models for compliance with the GDPR’s data minimization principle. This principle restricts prediction models to use the minimal information that is necessary for performing the task at hand. Given the challenge of the black-box setting, our key idea is to check if each of the prediction model’s input features is individually necessary by assigning it some constant value (i.e., applying a simple imputation) across all prediction instances, and measuring the extent to which the model outcomes would change. We introduce a metric for data minimization that is based on model instability under simple imputations. We extend the applicability of this metric from a finite sample model to a distributional setting by introducing a probabilistic data minimization guarantee, which we derive using a Bayesian approach. Furthermore, we address the auditing problem under a constraint on the number of queries to the prediction system. We formulate the problem of allocating a budget of system queries to feasible simple imputations (for investigating model instability) as a multi-armed bandit framework with probabilistic success metrics. We define two bandit problems for providing a probabilistic data minimization guarantee at a given confidence level: a decision problem given a data minimization level, and a measurement problem given a fixed query budget. We design efficient algorithms for these auditing problems using novel exploration strategies that expand classical bandit strategies. Our experiments with real-world prediction systems show that our auditing algorithms significantly outperform simpler benchmarks in both measurement and decision problems.

[Dueling Bandits with Team Comparisons](#)

- Lee Cohen, Ulrike Schmidt-Kraepelin, Yishay Mansour
- abstract@[open-review](#): We introduce the dueling teams problem, a new online-learning setting in which the learner observes noisy comparisons of disjoint pairs of \$k\$-sized teams from a universe of \$n\$ players. The goal of the learner is to minimize the number of duels required to identify, with high probability, a Condorcet winning team, i.e., a team which wins against any other disjoint team (with probability at least \$1/2\$).Noisy comparisons are linked to a total order on the teams. We formalize our model by building upon the dueling bandits setting (Yue et al. 2012) and provide several algorithms, both for stochastic and deterministic settings. For the stochastic setting, we provide a reduction to the classical dueling bandits setting, yielding an algorithm that identifies a Condorcet winning team within \$\\mathcal{O}((n+k \\log(k)) \\frac{\\max(\\log \\log n, \\log k)}{\\Delta^2})\$ duels, where \$\\Delta\$ is a gap parameter. For deterministic feedback, we additionally present a gap-independent algorithm that identifies a Condorcet winning team within \$\\mathcal{O}(nk \\log(k) + k^5)\$ duels.

[Meta Internal Learning](#)

- Raphael Bensadoun, Shir Gur, Tomer Galanti, Lior Wolf
- abstract@[open-review](#): Internal learning for single-image generation is a framework, where a generator is trained to produce novel images based on a single image. Since these models are trained on a single image, they are limited in their scale and application. To overcome these issues, we propose a meta-learning approach that enables training over a collection of images, in order to model the internal statistics of the sample image more effectively.In the presented meta-learning approach, a single-image GAN model is generated given an input image, via a convolutional feedforward hypernetwork \$f\$. This network is trained over a dataset of images, allowing for feature sharing among different models, and for interpolation in the space of generative models. The generated single-image model contains a hierarchy of multiple generators and discriminators. It is therefore required to train the meta-learner in an adversarial manner, which requires careful design choices that we justify by a theoretical analysis. Our results show that the models obtained are as

suitable as single-image GANs for many common image applications, {significantly reduce the training time per image without loss in performance}, and introduce novel capabilities, such as interpolation and feedforward modeling of novel images.

[Uniform Convergence of Interpolators: Gaussian Width, Norm Bounds and Benign Overfitting](#)

- Frederic Koehler, Lijia Zhou, Danica J. Sutherland, Nathan Srebro
- abstract@[open-review](#): We consider interpolation learning in high-dimensional linear regression with Gaussian data, and prove a generic uniform convergence guarantee on the generalization error of interpolators in an arbitrary hypothesis class in terms of the class's Gaussian width. Applying the generic bound to Euclidean norm balls recovers the consistency result of Bartlett et al. (2020) for minimum-norm interpolators, and confirms a prediction of Zhou et al. (2020) for near-minimal-norm interpolators in the special case of Gaussian data. We demonstrate the generality of the bound by applying it to the simplex, obtaining a novel consistency result for minimum $\|\cdot\|_1$ -norm interpolators (basis pursuit). Our results show how norm-based generalization bounds can explain and be used to analyze benign overfitting, at least in some settings.

[Adaptive wavelet distillation from neural networks through interpretations](#)

- Wooseok Ha, Chandan Singh, Francois Lanusse, Srigokul Upadhyayula, Bin Yu
- abstract@[open-review](#): Recent deep-learning models have achieved impressive prediction performance, but often sacrifice interpretability and computational efficiency. Interpretability is crucial in many disciplines, such as science and medicine, where models must be carefully vetted or where interpretation is the goal itself. Moreover, interpretable models are concise and often yield computational efficiency. Here, we propose adaptive wavelet distillation (AWD), a method which aims to distill information from a trained neural network into a wavelet transform. Specifically, AWD penalizes feature attributions of a neural network in the wavelet domain to learn an effective multi-resolution wavelet transform. The resulting model is highly predictive, concise, computationally efficient, and has properties (such as a multi-scale structure) which make it easy to interpret. In close collaboration with domain experts, we showcase how AWD addresses challenges in two real-world settings: cosmological parameter inference and molecular-partner prediction. In both cases, AWD yields a scientifically interpretable and concise model which gives predictive performance better than state-of-the-art neural networks. Moreover, AWD identifies predictive features that are scientifically meaningful in the context of respective domains. All code and models are released in a full-fledged package available on Github.

[Generative Occupancy Fields for 3D Surface-Aware Image Synthesis](#)

- Xudong XU, Xingang Pan, Dahua Lin, Bo Dai
- abstract@[open-review](#): The advent of generative radiance fields has significantly promoted the development of 3D-aware image synthesis. The cumulative rendering process in radiance fields makes training these generative models much easier since gradients are distributed over the entire volume, but leads to diffused object surfaces. In the meantime, compared to radiance fields occupancy representations could inherently ensure deterministic surfaces. However, if we directly apply occupancy representations to generative models, during training they will only receive sparse gradients located on object surfaces and eventually suffer from the convergence problem. In this paper, we propose Generative Occupancy Fields (GOF), a novel model based on generative radiance fields that can learn compact object surfaces without impeding its training convergence. The key insight of GOF is a dedicated transition from the cumulative rendering in radiance fields to rendering with only the surface points as the learned surface gets more and more accurate. In this way, GOF combines the merits of two representations in a unified framework. In practice, the training-time transition of start from radiance fields and march to occupancy representations is achieved in GOF by gradually shrinking the sampling region in its rendering process from the entire volume to a minimal neighboring region around the surface. Through comprehensive experiments on multiple datasets, we demonstrate that GOF can synthesize high-quality images with 3D consistency and simultaneously learn compact and smooth object surfaces. Our code is available at https://github.com/SheldonTsui/GOF_NeurIPS2021.

[Relaxed Marginal Consistency for Differentially Private Query Answering](#)

- Ryan McKenna, Siddhant Pradhan, Daniel R. Sheldon, Gerome Miklau
- abstract@[open-review](#): Many differentially private algorithms for answering database queries involve a step that reconstructs a discrete data distribution from noisy measurements. This provides consistent query answers and reduces error, but often requires space that grows exponentially with dimension. PRIVATE-PGM is a recent approach that uses graphical models to represent the data distribution, with complexity proportional to that of exact marginal inference in a graphical model with structure determined by the co-occurrence of variables in the noisy measurements. PRIVATE-PGM is highly scalable for sparse measurements, but may fail to run in high dimensions with dense measurements. We overcome the main scalability limitation of PRIVATE-PGM through a principled approach that relaxes consistency constraints in the estimation objective. Our new approach works with many existing private query answering algorithms and improves scalability or accuracy with no privacy cost.

[Local policy search with Bayesian optimization](#)

- Sarah Müller, Alexander von Rohr, Sebastian Trimpe
- abstract@[open-review](#): Reinforcement learning (RL) aims to find an optimal policy by interaction with an environment. Consequently, learning complex behavior requires a vast number of samples, which can be prohibitive in practice. Nevertheless, instead of systematically reasoning and actively choosing informative samples, policy gradients for local search are often obtained from random perturbations. These random samples yield high variance estimates and hence are sub-optimal in terms of sample complexity. Actively selecting informative samples is at the core of Bayesian optimization, which constructs a probabilistic surrogate of the objective from past samples to reason about informative subsequent ones. In this paper, we propose to join both worlds. We develop an algorithm utilizing a probabilistic model of the objective function and its gradient. Based on the model, the algorithm decides where to query a noisy zeroth-order oracle to improve the gradient estimates. The resulting algorithm is a novel type of policy search method, which we compare to existing black-box algorithms. The comparison reveals improved sample complexity and reduced variance in extensive empirical evaluations on synthetic objectives. Further, we highlight the benefits of active sampling on popular RL benchmarks.

[DominoSearch: Find layer-wise fine-grained N:M sparse schemes from dense neural networks](#)

- Wei Sun, Aojun Zhou, Sander Stuijk, Rob Wijhoven, Andrew Oakleigh Nelson, hongsheng Li, Henk Corporaal
- abstract@[open-review](#): Neural pruning is a widely-used compression technique for Deep Neural Networks (DNNs). Recent innovations in Hardware Architectures (e.g. Nvidia Ampere Sparse Tensor Core) and N:M fine-grained Sparse Neural Network algorithms (i.e. every M-weights contains N non-zero values) reveal a promising research line of neural pruning. However, the existing N:M algorithms only address the challenge of how to train N:M sparse neural networks in a uniform fashion (i.e. every layer has the same N:M sparsity) and suffer from a significant accuracy drop for high sparsity (i.e. when sparsity > 80%). To tackle this problem, we present a novel technique -- \textbf{\textit{DominoSearch}} to find mixed N:M sparsity schemes from pre-trained dense deep neural networks to achieve higher accuracy than the uniform-sparsity scheme with equivalent complexity constraints (e.g. model size or FLOPs). For instance, for the same model size with 2.1M parameters (87.5% sparsity), our layer-wise N:M sparse ResNet18 outperforms its uniform counterpart by 2.1% top-1 accuracy, on the large-scale ImageNet dataset. For the same computational complexity of 227M FLOPs, our layer-wise sparse ResNet18 outperforms the uniform one by 1.3% top-1 accuracy. Furthermore, our layer-wise fine-grained N:M sparse ResNet50 achieves 76.7% top-1 accuracy with 5.0M parameters. {This is competitive to the results achieved by layer-wise unstructured sparsity} that is believed to be the upper-bound of Neural Network pruning with respect to the accuracy-sparsity trade-off. We believe that our work can build a strong baseline for further

sparse DNN research and encourage future hardware-algorithm co-design work. Our code and models are publicly available at \url{https://github.com/NM-sparsity/DominoSearch}.

[Techniques for Symbol Grounding with SATNet](#)

- Sever Topan, David Rolnick, Xujie Si
- abstract@[open-review](#): Many experts argue that the future of artificial intelligence is limited by the field's ability to integrate symbolic logical reasoning into deep learning architectures. The recently proposed differentiable MAXSAT solver, SATNet, was a breakthrough in its capacity to integrate with a traditional neural network and solve visual reasoning problems. For instance, it can learn the rules of Sudoku purely from image examples. Despite its success, SATNet was shown to succumb to a key challenge in neurosymbolic systems known as the Symbol Grounding Problem: the inability to map visual inputs to symbolic variables without explicit supervision ("label leakage"). In this work, we present a self-supervised pre-training pipeline that enables SATNet to overcome this limitation, thus broadening the class of problems that SATNet architectures can solve to include datasets where no intermediary labels are available at all. We demonstrate that our method allows SATNet to attain full accuracy even with a harder problem setup that prevents any label leakage. We additionally introduce a proofreading method that further improves the performance of SATNet architectures, beating the state-of-the-art on Visual Sudoku.

[Object DGCNN: 3D Object Detection using Dynamic Graphs](#)

- Yue Wang, Justin M. Solomon
- abstract@[open-review](#): 3D object detection often involves complicated training and testing pipelines, which require substantial domain knowledge about individual datasets. Inspired by recent non-maximum suppression-free 2D object detection models, we propose a 3D object detection architecture on point clouds. Our method models 3D object detection as message passing on a dynamic graph, generalizing the DGCNN framework to predict a set of objects. In our construction, we remove the necessity of post-processing via object confidence aggregation or non-maximum suppression. To facilitate object detection from sparse point clouds, we also propose a set-to-set distillation approach customized to 3D detection. This approach aligns the outputs of the teacher model and the student model in a permutation-invariant fashion, significantly simplifying knowledge distillation for the 3D detection task. Our method achieves state-of-the-art performance on autonomous driving benchmarks. We also provide abundant analysis of the detection model and distillation framework.

[Safe Policy Optimization with Local Generalized Linear Function Approximations](#)

- Akifumi Wachi, Yunyue Wei, Yanan Sui
- abstract@[open-review](#): Safe exploration is a key to applying reinforcement learning (RL) in safety-critical systems. Existing safe exploration methods guaranteed safety under the assumption of regularity, and it has been difficult to apply them to large-scale real problems. We propose a novel algorithm, SPO-LF, that optimizes an agent's policy while learning the relation between a locally available feature obtained by sensors and environmental reward/safety using generalized linear function approximations. We provide theoretical guarantees on its safety and optimality. We experimentally show that our algorithm is 1) more efficient in terms of sample complexity and computational cost and 2) more applicable to large-scale problems than previous safe RL methods with theoretical guarantees, and 3) comparably sample-efficient and safer compared with existing advanced deep RL methods with safety constraints.

[Symplectic Adjoint Method for Exact Gradient of Neural ODE with Minimal Memory](#)

- Takashi Matsubara, Yuto Miyatake, Takaharu Yaguchi
- abstract@[open-review](#): A neural network model of a differential equation, namely neural ODE, has enabled the learning of continuous-time dynamical systems and probabilistic distributions with high accuracy. The neural ODE uses the same network repeatedly during a numerical integration. The memory consumption of the backpropagation algorithm is proportional to the number of uses times the network size. This is true even if a checkpointing scheme divides the computation graph into sub-graphs. Otherwise, the adjoint method obtains a gradient by a numerical integration backward in time. Although this method consumes memory only for a single network use, it requires high computational cost to suppress numerical errors. This study proposes the symplectic adjoint method, which is an adjoint method solved by a symplectic integrator. The symplectic adjoint method obtains the exact gradient (up to rounding error) with memory proportional to the number of uses plus the network size. The experimental results demonstrate that the symplectic adjoint method consumes much less memory than the naive backpropagation algorithm and checkpointing schemes, performs faster than the adjoint method, and is more robust to rounding errors.

[Exponential Separation between Two Learning Models and Adversarial Robustness](#)

- Grzegorz Gluch, Ruediger Urbanke
- abstract@[open-review](#): We prove an exponential separation for the sample/query complexity between the standard PAC-learning model and a version of the Equivalence-Query-learning model. In the PAC model all samples are provided at the beginning of the learning process. In the Equivalence-Query model the samples are acquired through an interaction between a teacher and a learner, where the teacher provides counterexamples to hypotheses given by the learner. It is intuitive that in an interactive setting fewer samples are needed. We make this formal and prove that in order to achieve an error $\$\\epsilon\$$ exponentially fewer samples suffice than what the PAC bound requires. It was shown experimentally by Stutz, Hein, and Schiele that adversarial training with on-manifold adversarial examples aids generalization (compared to standard training). If we think of the adversarial examples as counterexamples to the current hypothesis then our result can be thought of as a theoretical confirmation of those findings. We also discuss how our result relates to adversarial robustness. In the standard adversarial model one restricts the adversary by introducing a norm constraint. An alternative was pioneered by Goldwasser et. al. Rather than restricting the adversary the learner is enhanced. We pursue a third path. We require the adversary to return samples according to the Equivalence-Query model and show that this leads to robustness. Even though our model has its limitations it provides a fresh point of view on adversarial robustness.

[The balancing principle for parameter choice in distance-regularized domain adaptation](#)

- Werner Zellinger, Natalia Shepeleva, Marius-Constantin Dinu, Hamid Eghbal-zadeh, Hoan Duc Nguyen, Bernhard Nessler, Sergei Pereverzyev, Bernhard A. Moser
- abstract@[open-review](#): We address the unsolved algorithm design problem of choosing a justified regularization parameter in unsupervised domain adaptation. This problem is intriguing as no labels are available in the target domain. Our approach starts with the observation that the widely-used method of minimizing the source error, penalized by a distance measure between source and target feature representations, shares characteristics with regularized ill-posed inverse problems. Regularization parameters in inverse problems are optimally chosen by the fundamental principle of balancing approximation and sampling errors. We use this principle to balance learning errors and domain distance in a target error bound. As a result, we obtain a theoretically justified rule for the choice of the regularization parameter. In contrast to the state of the art, our approach allows source and target distributions with disjoint supports. An empirical comparative study on benchmark datasets underpins the performance of our approach.

[Gaussian Kernel Mixture Network for Single Image Defocus Deblurring](#)

- Yuhui Quan, Zicong Wu, Hui Ji
- abstract@[open-review](#): Defocus blur is one kind of blur effects often seen in images, which is challenging to remove due to its spatially variant amount. This paper presents an end-to-end deep learning approach for removing defocus blur from a single image, so as to have an all-in-focus image for consequent vision tasks. First, a pixel-wise Gaussian kernel mixture (GKM) model is proposed for representing spatially variant defocus blur kernels in an efficient linear parametric form, with higher accuracy than existing models. Then, a deep neural network called GKMNet is developed by unrolling a fixed-point iteration of the GKM-based deblurring. The GKMNet is built on a lightweight scale-recurrent architecture, with a scale-recurrent attention module for estimating the mixing coefficients in GKM for defocus deblurring. Extensive experiments show that the GKMNet not only noticeably outperforms existing defocus deblurring methods, but also has its advantages in terms of model complexity and computational efficiency.

[Cockpit: A Practical Debugging Tool for the Training of Deep Neural Networks](#)

- Frank Schneider, Felix Dangel, Philipp Hennig
- abstract@[open-review](#): When engineers train deep learning models, they are very much "flying blind". Commonly used methods for real-time training diagnostics, such as monitoring the train/test loss, are limited. Assessing a network's training process solely through these performance indicators is akin to debugging software without access to internal states through a debugger. To address this, we present Cockpit, a collection of instruments that enable a closer look into the inner workings of a learning machine, and a more informative and meaningful status report for practitioners. It facilitates the identification of learning phases and failure modes, like ill-chosen hyperparameters. These instruments leverage novel higher-order information about the gradient distribution and curvature, which has only recently become efficiently accessible. We believe that such a debugging tool, which we open-source for PyTorch, is a valuable help in troubleshooting the training process. By revealing new insights, it also more generally contributes to explainability and interpretability of deep nets.

[MEST: Accurate and Fast Memory-Economic Sparse Training Framework on the Edge](#)

- Geng Yuan, Xiaolong Ma, Wei Niu, Zhengang Li, Zhenglun Kong, Ning Liu, Yifan Gong, Zheng Zhan, Chaoyang He, Qing Jin, Siyue Wang, Minghai Qin, Bin Ren, Yanzhi Wang, Sijia Liu, Xue Lin
- abstract@[open-review](#): Recently, a new trend of exploring sparsity for accelerating neural network training has emerged, embracing the paradigm of training on the edge. This paper proposes a novel Memory-Economic Sparse Training (MEST) framework targeting for accurate and fast execution on edge devices. The proposed MEST framework consists of enhancements by Elastic Mutation (EM) and Soft Memory Bound (&S) that ensure superior accuracy at high sparsity ratios. Different from the existing works for sparse training, this current work reveals the importance of sparsity schemes on the performance of sparse training in terms of accuracy as well as training speed on real edge devices. On top of that, the paper proposes to employ data efficiency for further acceleration of sparse training. Our results suggest that unforgettable examples can be identified in-situ even during the dynamic exploration of sparsity masks in the sparse training process, and therefore can be removed for further training speedup on edge devices. Comparing with state-of-the-art (SOTA) works on accuracy, our MEST increases Top-1 accuracy significantly on ImageNet when using the same unstructured sparsity scheme. Systematical evaluation on accuracy, training speed, and memory footprint are conducted, where the proposed MEST framework consistently outperforms representative SOTA works. A reviewer strongly against our work based on his false assumptions and misunderstandings. On top of the previous submission, we employ data efficiency for further acceleration of sparse training. And we explore the impact of model sparsity, sparsity schemes, and sparse training algorithms on the number of removable training examples. Our codes are publicly available at: <https://github.com/boone891214/MEST>.

[Precise characterization of the prior predictive distribution of deep ReLU networks](#)

- Lorenzo Noci, Gregor Bachmann, Kevin Roth, Sebastian Nowozin, Thomas Hofmann
- abstract@[open-review](#): Recent works on Bayesian neural networks (BNNs) have highlighted the need to better understand the implications of using Gaussian priors in combination with the compositional structure of the network architecture. Similar in spirit to the kind of analysis that has been developed to devise better initialization schemes for neural networks (cf. He- or Xavier initialization), we derive a precise characterization of the prior predictive distribution of finite-width ReLU networks with Gaussian weights. While theoretical results have been obtained for their heavy-tailedness, the full characterization of the prior predictive distribution (i.e. its density, CDF and moments), remained unknown prior to this work. Our analysis, based on the Meijer-G function, allows us to quantify the influence of architectural choices such as the width or depth of the network on the resulting shape of the prior predictive distribution. We also formally connect our results to previous work in the infinite width setting, demonstrating that the moments of the distribution converge to those of a normal log-normal mixture in the infinite depth limit. Finally, our results provide valuable guidance on prior design: for instance, controlling the predictive variance with depth- and width-informed priors on the weights of the network.

[RED : Looking for Redundancies for Data-Free Structured Compression of Deep Neural Networks](#)

- Edouard YVINEC, Arnaud Dapogny, Matthieu Cord, Kevin Bailly
- abstract@[open-review](#): Deep Neural Networks (DNNs) are ubiquitous in today's computer vision landscape, despite involving considerable computational costs. The mainstream approaches for runtime acceleration consist in pruning connections (unstructured pruning) or, better, filters (structured pruning), both often requiring data to retrain the model. In this paper, we present RED, a data-free, unified approach to tackle structured pruning. First, we propose a novel adaptive hashing of the scalar DNN weight distribution densities to increase the number of identical neurons represented by their weight vectors. Second, we prune the network by merging redundant neurons based on their relative similarities, as defined by their distance. Third, we propose a novel uneven depthwise separation technique to further prune convolutional layers. We demonstrate through a large variety of benchmarks that RED largely outperforms other data-free pruning methods, often reaching performance similar to unconstrained, data-driven methods.

[TestRank: Bringing Order into Unlabeled Test Instances for Deep Learning Tasks](#)

- YU LI, Min LI, Qiuxia LAI, Yannan Liu, Qiang Xu
- abstract@[open-review](#): Deep learning (DL) systems are notoriously difficult to test and debug due to the lack of correctness proof and the huge test input space to cover. Given the ubiquitous unlabeled test data and high labeling cost, in this paper, we propose a novel test prioritization technique, namely TestRank, which aims at revealing more model failures with less labeling effort. TestRank brings order into the unlabeled test data according to their likelihood of being a failure, i.e., their failure-revealing capabilities. Different from existing solutions, TestRank leverages both intrinsic and contextual attributes of the unlabeled test data when prioritizing them. To be specific, we first build a similarity graph on both unlabeled test samples and labeled samples (e.g., training or previously labeled test samples). Then, we conduct graph-based semi-supervised learning to extract contextual features from the correctness of similar labeled samples. For a particular test instance, the contextual features extracted with the graph neural network and the intrinsic features obtained with the DL model itself are combined to predict its failure-revealing capability. Finally, TestRank prioritizes unlabeled test inputs in descending order of the above probability value. We evaluate TestRank on three popular image classification datasets, and results show that TestRank significantly outperforms existing test prioritization techniques.

[Large Scale Learning on Non-Homophilous Graphs: New Benchmarks and Strong Simple Methods](#)

- Derek Lim, Felix Hohne, Xiuyu Li, Sijia Linda Huang, Vaishnavi Gupta, Omkar Bhalerao, Ser Nam Lim
- abstract@[open-review](#): Many widely used datasets for graph machine learning tasks have generally been homophilous, where nodes with similar labels connect to each other. Recently, new Graph Neural Networks (GNNs) have been developed that move beyond the homophily regime; however, their

evaluation has often been conducted on small graphs with limited application domains. We collect and introduce diverse non-homophilous datasets from a variety of application areas that have up to 384x more nodes and 1398x more edges than prior datasets. We further show that existing scalable graph learning and graph minibatching techniques lead to performance degradation on these non-homophilous datasets, thus highlighting the need for further work on scalable non-homophilous methods. To address these concerns, we introduce LINKX --- a strong simple method that admits straightforward minibatch training and inference. Extensive experimental results with representative simple methods and GNNs across our proposed datasets show that LINKX achieves state-of-the-art performance for learning on non-homophilous graphs. Our codes and data are available at <https://github.com/CUAI/Non-Homophily-Large-Scale>.

[Reinforcement Learning based Disease Progression Model for Alzheimer's Disease](#)

- Krishnakant Saboo, Anirudh Choudhary, Yurui Cao, Gregory Worrell, David Jones, Ravishankar Iyer
- abstract@[open-review](#): We model Alzheimer's disease (AD) progression by combining differential equations (DEs) and reinforcement learning (RL) with domain knowledge. DEs provide relationships between some, but not all, factors relevant to AD. We assume that the missing relationships must satisfy general criteria about the working of the brain, for e.g., maximizing cognition while minimizing the cost of supporting cognition. This allows us to extract the missing relationships by using RL to optimize an objective (reward) function that captures the above criteria. We use our model consisting of DEs (as a simulator) and the trained RL agent to predict individualized 10-year AD progression using baseline (year 0) features on synthetic and real data. The model was comparable or better at predicting 10-year cognition trajectories than state-of-the-art learning-based models. Our interpretable model demonstrated, and provided insights into, "recovery/compensatory" processes that mitigate the effect of AD, even though those processes were not explicitly encoded in the model. Our framework combines DEs with RL for modelling AD progression and has broad applicability for understanding other neurological disorders.

[Catch-A-Waveform: Learning to Generate Audio from a Single Short Example](#)

- Gal Greshler, Tamar Shaham, Tomer Michaeli
- abstract@[open-review](#): Models for audio generation are typically trained on hours of recordings. Here, we illustrate that capturing the essence of an audio source is typically possible from as little as a few tens of seconds from a single training signal. Specifically, we present a GAN-based generative model that can be trained on one short audio signal from any domain (e.g. speech, music, etc.) and does not require pre-training or any other form of external supervision. Once trained, our model can generate random samples of arbitrary duration that maintain semantic similarity to the training waveform, yet exhibit new compositions of its audio primitives. This enables a long line of interesting applications, including generating new jazz improvisations or new a-cappella rap variants based on a single short example, producing coherent modifications to famous songs (e.g. adding a new verse to a Beatles song based solely on the original recording), filling-in of missing parts (inpainting), extending the bandwidth of a speech signal (super-resolution), and enhancing old recordings without access to any clean training example. We show that in all cases, no more than 20 seconds of training audio commonly suffice for our model to achieve state-of-the-art results. This is despite its complete lack of prior knowledge about the nature of audio signals in general.

[Explanation-based Data Augmentation for Image Classification](#)

- Sandareka Wickramanayake, Wynne Hsu, Mong Li Lee
- abstract@[open-review](#): Existing works have generated explanations for deep neural network decisions to provide insights into model behavior. We observe that these explanations can also be used to identify concepts that caused misclassifications. This allows us to understand the possible limitations of the dataset used to train the model, particularly the under-represented regions in the dataset. This work proposes a framework that utilizes concept-based explanations to automatically augment the dataset with new images that can cover these under-represented regions to improve the model performance. The framework is able to use the explanations generated by both interpretable classifiers and post-hoc explanations from black-box classifiers. Experiment results demonstrate that the proposed approach improves the accuracy of classifiers compared to state-of-the-art augmentation strategies.

[Data-Efficient GAN Training Beyond \(Just\) Augmentations: A Lottery Ticket Perspective](#)

- Tianlong Chen, Yu Cheng, Zhe Gan, Jingjing Liu, Zhangyang Wang
- abstract@[open-review](#): Training generative adversarial networks (GANs) with limited real image data generally results in deteriorated performance and collapsed models. To conquer this challenge, we are inspired by the latest observation, that one can discover independently trainable and highly sparse subnetworks (a.k.a., lottery tickets) from GANs. Treating this as an inductive prior, we suggest a brand-new angle towards data-efficient GAN training: by first identifying the lottery ticket from the original GAN using the small training set of real images; and then focusing on training that sparse subnetwork by re-using the same set. We find our coordinated framework to offer orthogonal gains to existing real image data augmentation methods, and we additionally present a new feature-level augmentation that can be applied together with them. Comprehensive experiments endorse the effectiveness of our proposed framework, across various GAN architectures (SNGAN, BigGAN, and StyleGAN-V2) and diverse datasets (CIFAR-10, CIFAR-100, Tiny-ImageNet, ImageNet, and multiple few-shot generation datasets). Codes are available at: <https://github.com/VITA-Group/Ultra-Data-Efficient-GAN-Training>.

[When Are Solutions Connected in Deep Networks?](#)

- Quynh N. Nguyen, Pierre Bréchet, Marco Mondelli
- abstract@[open-review](#): The question of how and why the phenomenon of mode connectivity occurs in training deep neural networks has gained remarkable attention in the research community. From a theoretical perspective, two possible explanations have been proposed: (i) the loss function has connected sublevel sets, and (ii) the solutions found by stochastic gradient descent are dropout stable. While these explanations provide insights into the phenomenon, their assumptions are not always satisfied in practice. In particular, the first approach requires the network to have one layer with order of $\$N\$$ neurons ($\$N$$ being the number of training samples), while the second one requires the loss to be almost invariant after removing half of the neurons at each layer (up to some rescaling of the remaining ones). In this work, we improve both conditions by exploiting the quality of the features at every intermediate layer together with a milder over-parameterization requirement. More specifically, we show that: (i) under generic assumptions on the features of intermediate layers, it suffices that the last two hidden layers have order of $\$sqrt{N}\$$ neurons, and (ii) if subsets of features at each layer are linearly separable, then almost no over-parameterization is needed to show the connectivity. Our experiments confirm that the proposed condition ensures the connectivity of solutions found by stochastic gradient descent, even in settings where the previous requirements do not hold.

[TOHAN: A One-step Approach towards Few-shot Hypothesis Adaptation](#)

- Haoang Chi, Feng Liu, Wenjing Yang, Long Lan, Tongliang Liu, Bo Han, William Cheung, James Kwok
- abstract@[open-review](#): In few-shot domain adaptation (FDA), classifiers for the target domain are trained with \emph{accessible} labeled data in the source domain (SD) and few labeled data in the target domain (TD). However, data usually contain private information in the current era, e.g., data distributed on personal phones. Thus, the private data will be leaked if we directly access data in SD to train a target-domain classifier (required by FDA methods). In this paper, to prevent privacy leakage in SD, we consider a very challenging problem setting, where the classifier for the TD has to be trained using few labeled target data and a well-trained SD classifier, named few-shot hypothesis adaptation (FHA). In FHA, we cannot access data in SD, as a result, the private information in SD will be protected well. To this end, we propose a target-oriented hypothesis adaptation network (TOHAN) to solve the FHA problem, where we generate highly-compatible unlabeled data (i.e., an intermediate domain) to help train a target-domain classifier. TOHAN

maintains two deep networks simultaneously, in which one focuses on learning an intermediate domain and the other takes care of the intermediate-to-target distributional adaptation and the target-risk minimization. Experimental results show that TOHAN outperforms competitive baselines significantly.

[Learning Graph Cellular Automata](#)

- Daniele Grattarola, Lorenzo Livi, Cesare Alippi
- abstract@[open-review](#): Cellular automata (CA) are a class of computational models that exhibit rich dynamics emerging from the local interaction of cells arranged in a regular lattice. In this work we focus on a generalised version of typical CA, called graph cellular automata (GCA), in which the lattice structure is replaced by an arbitrary graph. In particular, we extend previous work that used convolutional neural networks to learn the transition rule of conventional CA and we use graph neural networks to learn a variety of transition rules for GCA. First, we present a general-purpose architecture for learning GCA, and we show that it can represent any arbitrary GCA with finite and discrete state space. Then, we test our approach on three different tasks: 1) learning the transition rule of a GCA on a Voronoi tessellation; 2) imitating the behaviour of a group of flocking agents; 3) learning a rule that converges to a desired target state.

[Efficient Online Estimation of Causal Effects by Deciding What to Observe](#)

- Shantanu Gupta, Zachary Lipton, David Childers
- abstract@[open-review](#): Researchers often face data fusion problems, where multiple data sources are available, each capturing a distinct subset of variables. While problem formulations typically take the data as given, in practice, data acquisition can be an ongoing process. In this paper, we introduce the problem of deciding, at each time, which data source to sample from. Our goal is to estimate a given functional of the parameters of a probabilistic model as efficiently as possible. We propose online moment selection (OMS), a framework in which structural assumptions are encoded as moment conditions. The optimal action at each step depends, in part, on the very moments that identify the functional of interest. Our algorithms balance exploration with choosing the best action as suggested by estimated moments. We propose two selection strategies: (1) explore-then-commit (ETC) and (2) explore-then-greedy (ETG), proving that both achieve zero asymptotic regret as assessed by MSE. We instantiate our setup for average treatment effect estimation, where structural assumptions are given by a causal graph and data sources include subsets of mediators, confounders, and instrumental variables.

[Perturbation Theory for the Information Bottleneck](#)

- Vudtiwat Ngampruetikorn, David J. Schwab
- abstract@[open-review](#): Extracting relevant information from data is crucial for all forms of learning. The information bottleneck (IB) method formalizes this, offering a mathematically precise and conceptually appealing framework for understanding learning phenomena. However the nonlinearity of the IB problem makes it computationally expensive and analytically intractable in general. Here we derive a perturbation theory for the IB method and report the first complete characterization of the learning onset, the limit of maximum relevant information per bit extracted from data. We test our results on synthetic probability distributions, finding good agreement with the exact numerical solution near the onset of learning. We explore the difference and subtleties in our derivation and previous attempts at deriving a perturbation theory for the learning onset and attribute the discrepancy to a flawed assumption. Our work also provides a fresh perspective on the intimate relationship between the IB method and the strong data processing inequality.

[Deconvolutional Networks on Graph Data](#)

- Jia Li, Jiajin Li, Yang Liu, Jianwei Yu, Yueling Li, Hong Cheng
- abstract@[open-review](#): In this paper, we consider an inverse problem in graph learning domain -- "given the graph representations smoothed by Graph Convolutional Network (GCN), how can we reconstruct the input graph signal?" We propose Graph Deconvolutional Network (GDN) and motivate the design of GDN via a combination of inverse filters in spectral domain and de-noising layers in wavelet domain, as the inverse operation results in a high frequency amplifier and may amplify the noise. We demonstrate the effectiveness of the proposed method on several tasks including graph feature imputation and graph structure generation.

[Variational Multi-Task Learning with Gumbel-Softmax Priors](#)

- Jiayi Shen, Xiantong Zhen, Marcel Worring, Ling Shao
- abstract@[open-review](#): Multi-task learning aims to explore task relatedness to improve individual tasks, which is of particular significance in the challenging scenario that only limited data is available for each task. To tackle this challenge, we propose variational multi-task learning (VMTL), a general probabilistic inference framework for learning multiple related tasks. We cast multi-task learning as a variational Bayesian inference problem, in which task relatedness is explored in a unified manner by specifying priors. To incorporate shared knowledge into each task, we design the prior of a task to be a learnable mixture of the variational posteriors of other related tasks, which is learned by the Gumbel-Softmax technique. In contrast to previous methods, our VMTL can exploit task relatedness for both representations and classifiers in a principled way by jointly inferring their posteriors. This enables individual tasks to fully leverage inductive biases provided by related tasks, therefore improving the overall performance of all tasks. Experimental results demonstrate that the proposed VMTL is able to effectively tackle a variety of challenging multi-task learning settings with limited training data for both classification and regression. Our method consistently surpasses previous methods, including strong Bayesian approaches, and achieves state-of-the-art performance on five benchmark datasets.

[Accelerating Quadratic Optimization with Reinforcement Learning](#)

- Jeffrey Ichniowski, Paras Jain, Bartolomeo Stellato, Goran Banjac, Michael Luo, Francesco Borrelli, Joseph E. Gonzalez, Ion Stoica, Ken Goldberg
- abstract@[open-review](#): First-order methods for quadratic optimization such as OSQP are widely used for large-scale machine learning and embedded optimal control, where many related problems must be rapidly solved. These methods face two persistent challenges: manual hyperparameter tuning and convergence time to high-accuracy solutions. To address these, we explore how Reinforcement Learning (RL) can learn a policy to tune parameters to accelerate convergence. In experiments with well-known QP benchmarks we find that our RL policy, RLQP, significantly outperforms state-of-the-art QP solvers by up to 3x. RLQP generalizes surprisingly well to previously unseen problems with varying dimension and structure from different applications, including the QPLIB, Netlib LP and Maros-M{\'e}sz{\'a}ros problems. Code, models, and videos are available at <https://berkeleyautomation.github.io/rlqp/>.

[Deep Residual Learning in Spiking Neural Networks](#)

- Wei Fang, Zhaofei Yu, Yanqi Chen, Tiejun Huang, Timoth{\'e}e Masquelier, Yonghong Tian
- abstract@[open-review](#): Deep Spiking Neural Networks (SNNs) present optimization difficulties for gradient-based approaches due to discrete binary activation and complex spatial-temporal dynamics. Considering the huge success of ResNet in deep learning, it would be natural to train deep SNNs with residual learning. Previous Spiking ResNet mimics the standard residual block in ANNs and simply replaces ReLU activation layers with spiking neurons, which suffers the degradation problem and can hardly implement residual learning. In this paper, we propose the spike-element-wise (SEW) ResNet to realize residual learning in deep SNNs. We prove that the SEW ResNet can easily implement identity mapping and overcome the vanishing/exploding gradient problems of Spiking ResNet. We evaluate our SEW ResNet on ImageNet, DVS Gesture, and CIFAR10-DVS datasets, and show that SEW

ResNet outperforms the state-of-the-art directly trained SNNs in both accuracy and time-steps. Moreover, SEW ResNet can achieve higher performance by simply adding more layers, providing a simple method to train deep SNNs. To our best knowledge, this is the first time that directly training deep SNNs with more than 100 layers becomes possible. Our codes are available at <https://github.com/fangwei123456/Spike-Element-Wise-ResNet>.

[Duplex Sequence-to-Sequence Learning for Reversible Machine Translation](#)

- Zaixiang Zheng, Hao Zhou, Shujian Huang, Jiajun Chen, Jingjing Xu, Lei Li
- abstract@[open-review](#): Sequence-to-sequence learning naturally has two directions. How to effectively utilize supervision signals from both directions? Existing approaches either require two separate models, or a multitask-learned model but with inferior performance. In this paper, we propose REDER (Reversible Duplex Transformer), a parameter-efficient model and apply it to machine translation. Either end of REDER can simultaneously input and output a distinct language. Thus REDER enables {\em reversible machine translation} by simply flipping the input and output ends. Experiments verify that REDER achieves the first success of reversible machine translation, which helps outperform its multitask-trained baselines by up to 1.3 BLEU.

[Improved Coresets and Sublinear Algorithms for Power Means in Euclidean Spaces](#)

- Vincent Cohen-Addad, David Saulpic, Chris Schwiegelshohn
- abstract@[open-review](#): In this paper, we consider the problem of finding high dimensional power means: given a set A of n points in \mathbb{R}^d , find the point m that minimizes the sum of Euclidean distance, raised to the power z , over all input points. Special cases of problem include the well-known Fermat-Weber problem -- or geometric median problem -- where $z = 1$, the mean or centroid where $z=2$, and the Minimum Enclosing Ball problem, where $z = \infty$. We consider these problem in the big data regime. Here, we are interested in sampling as few points as possible such that we can accurately estimate m . More specifically, we consider sublinear algorithms as well as coresets for these problems. Sublinear algorithms have a random query access to the A and the goal is to minimize the number of queries. Here, we show that $\tilde{O}(\epsilon^{-z-3})$ samples are sufficient to achieve a $(1+\epsilon)$ approximation, generalizing the results from Cohen, Lee, Miller, Pachocki, and Sidford [STOC '16] and Inaba, Katoh, and Imai [SoCG '94] to arbitrary z . Moreover, we show that this bound is nearly optimal, as any algorithm requires at least $\Omega(\epsilon^{-z+1})$ queries to achieve said approximation. The second contribution are coresets for these problems, where we aim to find a small, weighted subset of the points which approximate cost of every candidate point $c \in \mathbb{R}^d$ up to a $(1+\epsilon)$ factor. Here, we show that $\tilde{O}(\epsilon^{-2})$ points are sufficient, improving on the $\tilde{O}(d\epsilon^{-2})$ bound by Feldman and Langberg [STOC '11] and the $\tilde{O}(\epsilon^{-4})$ bound by Braverman, Jiang, Krauthamer, and Wu [SODA '21].

[Accelerated Sparse Neural Training: A Provable and Efficient Method to Find N:M Transposable Masks](#)

- Itay Hubara, Brian Chmiel, Moshe Island, Ron Banner, Joseph Naor, Daniel Soudry
- abstract@[open-review](#): Unstructured pruning reduces the memory footprint in deep neural networks (DNNs). Recently, researchers proposed different types of structural pruning intending to reduce also the computation complexity. In this work, we first suggest a new measure called mask-diversity which correlates with the expected accuracy of the different types of structural pruning. We focus on the recently suggested N:M fine-grained block sparsity mask, in which for each block of M weights, we have at least N zeros. While N:M fine-grained block sparsity allows acceleration in actual modern hardware, it can be used only to accelerate the inference phase. In order to allow for similar accelerations in the training phase, we suggest a novel transposable fine-grained sparsity mask, where the same mask can be used for both forward and backward passes. Our transposable mask guarantees that both the weight matrix and its transpose follow the same sparsity pattern; thus, the matrix multiplication required for passing the error backward can also be accelerated. We formulate the problem of finding the optimal transposable-mask as a minimum-cost flow problem. Additionally, to speed up the minimum-cost flow computation, we also introduce a fast linear-time approximation that can be used when the masks dynamically change during training. Our experiments suggest a 2x speed-up in the matrix multiplications with no accuracy degradation over vision and language models. Finally, to solve the problem of switching between different structure constraints, we suggest a method to convert a pre-trained model with unstructured sparsity to an N:M fine-grained block sparsity model with little to no training. A reference implementation can be found at <https://github.com/papers-submission/structuredtransposablemasks>.

[Learning and Generalization in RNNs](#)

- Abhishek Panigrahi, Navin Goyal
- abstract@[open-review](#): Simple recurrent neural networks (RNNs) and their more advanced cousins LSTMs etc. have been very successful in sequence modeling. Their theoretical understanding, however, is lacking and has not kept pace with the progress for feedforward networks, where a reasonably complete understanding in the special case of highly overparametrized one-hidden-layer networks has emerged. In this paper, we make progress towards remedying this situation by proving that RNNs can learn functions of sequences. In contrast to the previous work that could only deal with functions of sequences that are sums of functions of individual tokens in the sequence, we allow general functions. Conceptually and technically, we introduce new ideas which enable us to extract information from the hidden state of the RNN in our proofs---addressing a crucial weakness in previous work. We illustrate our results on some regular language recognition problems.

[Improving Visual Quality of Image Synthesis by A Token-based Generator with Transformers](#)

- Yanhong Zeng, Huan Yang, Hongyang Chao, Jianbo Wang, Jianlong Fu
- abstract@[open-review](#): We present a new perspective of achieving image synthesis by viewing this task as a visual token generation problem. Different from existing paradigms that directly synthesize a full image from a single input (e.g., a latent code), the new formulation enables a flexible local manipulation for different image regions, which makes it possible to learn content-aware and fine-grained style control for image synthesis. Specifically, it takes as input a sequence of latent tokens to predict the visual tokens for synthesizing an image. Under this perspective, we propose a token-based generator (i.e., TokenGAN). Particularly, the TokenGAN inputs two semantically different visual tokens, i.e., the learned constant content tokens and the style tokens from the latent space. Given a sequence of style tokens, the TokenGAN is able to control the image synthesis by assigning the styles to the content tokens by attention mechanism with a Transformer. We conduct extensive experiments and show that the proposed TokenGAN has achieved state-of-the-art results on several widely-used image synthesis benchmarks, including FFHQ and LSUN CHURCH with different resolutions. In particular, the generator is able to synthesize high-fidelity images with (1024x1024) size, dispensing with convolutions entirely.

[The Effect of the Intrinsic Dimension on the Generalization of Quadratic Classifiers](#)

- Fabian Latorre, Leello Tadesse Dadi, Paul Rolland, Volkan Cevher
- abstract@[open-review](#): It has been recently observed that neural networks, unlike kernel methods, enjoy a reduced sample complexity when the distribution is isotropic (i.e., when the covariance matrix is the identity). We find that this sensitivity to the data distribution is not exclusive to neural networks, and the same phenomenon can be observed on the class of quadratic classifiers (i.e., the sign of a quadratic polynomial) with a nuclear-norm constraint. We demonstrate this by deriving an upper bound on the Rademacher Complexity that depends on two key quantities: (i) the intrinsic dimension, which is a measure of isotropy, and (ii) the largest eigenvalue of the second moment (covariance) matrix of the distribution. Our result improves the dependence on the dimension over the best previously known bound and precisely quantifies the relation between the sample complexity and the level of isotropy of the distribution.

[DeepReduce: A Sparse-tensor Communication Framework for Federated Deep Learning](#)

- Hang Xu, Kelly Kostopoulou, Aritra Dutta, Xin Li, Alexandros Ntoulas, Panos Kalnis
- abstract@[open-review](#): Sparse tensors appear frequently in federated deep learning, either as a direct artifact of the deep neural network's gradients, or as a result of an explicit sparsification process. Existing communication primitives are agnostic to the peculiarities of deep learning; consequently, they impose unnecessary communication overhead. This paper introduces DeepReduce, a versatile framework for the compressed communication of sparse tensors, tailored to federated deep learning. DeepReduce decomposes sparse tensors into two sets, values and indices, and allows both independent and combined compression of these sets. We support a variety of common compressors, such as Deflate for values, or run-length encoding for indices. We also propose two novel compression schemes that achieve superior results: curve fitting-based for values, and bloom filter-based for indices. DeepReduce is orthogonal to existing gradient sparsifiers and can be applied in conjunction with them, transparently to the end-user, to significantly lower the communication overhead. As proof of concept, we implement our approach on TensorFlow and PyTorch. Our experiments with large real models demonstrate that DeepReduce transmits 320% less data than existing sparsifiers, without affecting accuracy. Code is available at <https://github.com/hangxu0304/DeepReduce>.

[Provably Efficient Causal Reinforcement Learning with Confounded Observational Data](#)

- Lingxiao Wang, Zhuoran Yang, Zhaoran Wang
- abstract@[open-review](#): Empowered by neural networks, deep reinforcement learning (DRL) achieves tremendous empirical success. However, DRL requires a large dataset by interacting with the environment, which is unrealistic in critical scenarios such as autonomous driving and personalized medicine. In this paper, we study how to incorporate the dataset collected in the offline setting to improve the sample efficiency in the online setting. To incorporate the observational data, we face two challenges. (a) The behavior policy that generates the observational data may depend on unobserved random variables (confounders), which affect the received rewards and transition dynamics. (b) Exploration in the online setting requires quantifying the uncertainty given both the observational and interventional data. To tackle such challenges, we propose the deconfounded optimistic value iteration (DOVI) algorithm, which incorporates the confounded observational data in a provably efficient manner. DOVI explicitly adjusts for the confounding bias in the observational data, where the confounders are partially observed or unobserved. In both cases, such adjustments allow us to construct the bonus based on a notion of information gain, which takes into account the amount of information acquired from the offline setting. In particular, we prove that the regret of DOVI is smaller than the optimal regret achievable in the pure online setting when the confounded observational data are informative upon the adjustments.

[Predicting Deep Neural Network Generalization with Perturbation Response Curves](#)

- Yair Schiff, Brian Quanz, Payel Das, Pin-Yu Chen
- abstract@[open-review](#): The field of Deep Learning is rich with empirical evidence of human-like performance on a variety of prediction tasks. However, despite these successes, the recent Predicting Generalization in Deep Learning (PGDL) NeurIPS 2020 competition suggests that there is a need for more robust and efficient measures of network generalization. In this work, we propose a new framework for evaluating the generalization capabilities of trained networks. We use perturbation response (PR) curves that capture the accuracy change of a given network as a function of varying levels of training sample perturbation. From these PR curves, we derive novel statistics that capture generalization capability. Specifically, we introduce two new measures for accurately predicting generalization gaps: the Gi-score and Pal-score, which are inspired by the Gini coefficient and Palma ratio (measures of income inequality), that accurately predict generalization gaps. Using our framework applied to intra and inter-class sample mixup, we attain better predictive scores than the current state-of-the-art measures on a majority of tasks in the PGDL competition. In addition, we show that our framework and the proposed statistics can be used to capture to what extent a trained network is invariant to a given parametric input transformation, such as rotation or translation. Therefore, these generalization gap prediction statistics also provide a useful means for selecting optimal network architectures and hyperparameters that are invariant to a certain perturbation.

[Exploiting Domain-Specific Features to Enhance Domain Generalization](#)

- Manh-Ha Bui, Toan Tran, Anh Tran, Dinh Phung
- abstract@[open-review](#): Domain Generalization (DG) aims to train a model, from multiple observed source domains, in order to perform well on unseen target domains. To obtain the generalization capability, prior DG approaches have focused on extracting domain-invariant information across sources to generalize on target domains, while useful domain-specific information which strongly correlates with labels in individual domains and the generalization to target domains is usually ignored. In this paper, we propose meta-Domain Specific-Domain Invariant (mDSDI) - a novel theoretically sound framework that extends beyond the invariance view to further capture the usefulness of domain-specific information. Our key insight is to disentangle features in the latent space while jointly learning both domain-invariant and domain-specific features in a unified framework. The domain-specific representation is optimized through the meta-learning framework to adapt from source domains, targeting a robust generalization on unseen domains. We empirically show that mDSDI provides competitive results with state-of-the-art techniques in DG. A further ablation study with our generated dataset, Background-Colored-MNIST, confirms the hypothesis that domain-specific is essential, leading to better results when compared with only using domain-invariant.

[Optimal Order Simple Regret for Gaussian Process Bandits](#)

- Sattar Vakili, Nacime Bouziani, Sepehr Jalali, Alberto Bernacchia, Da-shan Shiu
- abstract@[open-review](#): Consider the sequential optimization of a continuous, possibly non-convex, and expensive to evaluate objective function f . The problem can be cast as a Gaussian Process (GP) bandit where f lives in a reproducing kernel Hilbert space (RKHS). The state of the art analysis of several learning algorithms shows a significant gap between the lower and upper bounds on the simple regret performance. When N is the number of exploration trials and γ_N is the maximal information gain, we prove an $\tilde{O}(\sqrt{\gamma_N})$ bound on the simple regret performance of a pure exploration algorithm that is significantly tighter than the existing bounds. We show that this bound is order optimal up to logarithmic factors for the cases where a lower bound on regret is known. To establish these results, we prove novel and sharp confidence intervals for GP models applicable to RKHS elements which may be of broader interest.

[Generalization Guarantee of SGD for Pairwise Learning](#)

- Yunwen Lei, Mingrui Liu, Yiming Ying
- abstract@[open-review](#): Recently, there is a growing interest in studying pairwise learning since it includes many important machine learning tasks as specific examples, e.g., metric learning, AUC maximization and ranking. While stochastic gradient descent (SGD) is an efficient method, there is a lacking study on its generalization behavior for pairwise learning. In this paper, we present a systematic study on the generalization analysis of SGD for pairwise learning to understand the balance between generalization and optimization. We develop a novel high-probability generalization bound for uniformly-stable algorithms to incorporate the variance information for better generalization, based on which we establish the first nonsmooth learning algorithm to achieve almost optimal high-probability and dimension-independent generalization bounds in linear time. We consider both convex and nonconvex pairwise learning problems. Our stability analysis for convex problems shows how the interpolation can help generalization. We establish a uniform convergence of gradients, and apply it to derive the first generalization bounds on population gradients for nonconvex problems. Finally, we develop better generalization bounds for gradient-dominated problems.

[Supercharging Imbalanced Data Learning With Energy-based Contrastive Representation Transfer](#)

- Junya Chen, Zidi Xiu, Benjamin Goldstein, Ricardo Henao, Lawrence Carin, Chenyang Tao
- abstract@[open-review](#): Dealing with severe class imbalance poses a major challenge for many real-world applications, especially when the accurate classification and generalization of minority classes are of primary interest. In computer vision and NLP, learning from datasets with long-tail behavior is a recurring theme, especially for naturally occurring labels. Existing solutions mostly appeal to sampling or weighting adjustments to alleviate the extreme imbalance, or impose inductive bias to prioritize generalizable associations. Here we take a novel perspective to promote sample efficiency and model generalization based on the invariance principles of causality. Our contribution posits a meta-distributional scenario, where the causal generating mechanism for label-conditional features is invariant across different labels. Such causal assumption enables efficient knowledge transfer from the dominant classes to their under-represented counterparts, even if their feature distributions show apparent disparities. This allows us to leverage a causal data augmentation procedure to enlarge the representation of minority classes. Our development is orthogonal to the existing imbalanced data learning techniques thus can be seamlessly integrated. The proposed approach is validated on an extensive set of synthetic and real-world tasks against state-of-the-art solutions.

Heavy Ball Momentum for Conditional Gradient

- Bingcong Li, Alireza Sadeghi, Georgios Giannakis
- abstract@[open-review](#): Conditional gradient, aka Frank Wolfe (FW) algorithms, have well-documented merits in machine learning and signal processing applications. Unlike projection-based methods, momentum cannot improve the convergence rate of FW, in general. This limitation motivates the present work, which deals with heavy ball momentum, and its impact to FW. Specifically, it is established that heavy ball offers a unifying perspective on the primal-dual (PD) convergence, and enjoys a tighter \textit{per iteration} PD error rate, for multiple choices of step sizes, where PD error can serve as the stopping criterion in practice. In addition, it is asserted that restart, a scheme typically employed jointly with Nesterov's momentum, can further tighten this PD error bound. Numerical results demonstrate the usefulness of heavy ball momentum in FW iterations.

PARP: Prune, Adjust and Re-Prune for Self-Supervised Speech Recognition

- Cheng-I Jeff Lai, Yang Zhang, Alexander H. Liu, Shiyu Chang, Yi-Lun Liao, Yung-Sung Chuang, Kaizhi Qian, Sameer Khurana, David Cox, Jim Glass
- abstract@[open-review](#): Self-supervised speech representation learning (speech SSL) has demonstrated the benefit of scale in learning rich representations for Automatic Speech Recognition (ASR) with limited paired data, such as wav2vec 2.0. We investigate the existence of sparse subnetworks in pre-trained speech SSL models that achieve even better low-resource ASR results. However, directly applying widely adopted pruning methods such as the Lottery Ticket Hypothesis (LTH) is suboptimal in the computational cost needed. Moreover, we show that the discovered subnetworks yield minimal performance gain compared to the original dense network. We present Prune-Adjust-Re-Prune (PARP), which discovers and finetunes subnetworks for much better performance, while only requiring a single downstream ASR finetuning run. PARP is inspired by our surprising observation that subnetworks pruned for pre-training tasks need merely a slight adjustment to achieve a sizeable performance boost in downstream ASR tasks. Extensive experiments on low-resource ASR verify (1) sparse subnetworks exist in mono-lingual/multi-lingual pre-trained speech SSL, and (2) the computational advantage and performance gain of PARP over baseline pruning methods. In particular, on the 10min Librispeech split without LM decoding, PARP discovers subnetworks from wav2vec 2.0 with an absolute 10.9%/12.6% WER decrease compared to the full model. We further demonstrate the effectiveness of PARP via: cross-lingual pruning without any phone recognition degradation, the discovery of a multi-lingual subnetwork for 10 spoken languages in 1 finetuning run, and its applicability to pre-trained BERT/XLNet for natural language tasks1.

Robust Learning of Optimal Auctions

- Wenshuo Guo, Michael Jordan, Emmanouil Zampetakis
- abstract@[open-review](#): We study the problem of learning revenue-optimal multi-bidder auctions from samples when the samples of bidders' valuations can be adversarially corrupted or drawn from distributions that are adversarially perturbed. First, we prove tight upper bounds on the revenue we can obtain with a corrupted distribution under a population model, for both regular valuation distributions and distributions with monotone hazard rate (MHR). We then propose new algorithms that, given only an approximate distribution '' for the bidder's valuation, can learn a mechanism whose revenue is nearly optimal simultaneously for all true distributions'' that are α -close to the original distribution in Kolmogorov-Smirnov distance. The proposed algorithms operate beyond the setting of bounded distributions that have been studied in prior works, and are guaranteed to obtain a fraction $1-O(\alpha)$ of the optimal revenue under the true distribution when the distributions are MHR. Moreover, they are guaranteed to yield at least a fraction $1-O(\sqrt{\alpha})$ of the optimal revenue when the distributions are regular. We prove that these upper bounds cannot be further improved, by providing matching lower bounds. Lastly, we derive sample complexity upper bounds for learning a near-optimal auction for both MHR and regular distributions.

Disrupting Deep Uncertainty Estimation Without Harming Accuracy

- Ido Galil, Ran El-Yaniv
- abstract@[open-review](#): Deep neural networks (DNNs) have proven to be powerful predictors and are widely used for various tasks. Credible uncertainty estimation of their predictions, however, is crucial for their deployment in many risk-sensitive applications. In this paper we present a novel and simple attack, which unlike adversarial attacks, does not cause incorrect predictions but instead cripples the network's capacity for uncertainty estimation. The result is that after the attack, the DNN is more confident of its incorrect predictions than about its correct ones without having its accuracy reduced. We present two versions of the attack. The first scenario focuses on a black-box regime (where the attacker has no knowledge of the target network) and the second scenario attacks a white-box setting. The proposed attack is only required to be of minuscule magnitude for its perturbations to cause severe uncertainty estimation damage, with larger magnitudes resulting in completely unusable uncertainty estimations. We demonstrate successful attacks on three of the most popular uncertainty estimation methods: the vanilla softmax score, Deep Ensembles and MC-Dropout. Additionally, we show an attack on SelectiveNet, the selective classification architecture. We test the proposed attack on several contemporary architectures such as MobileNetV2 and EfficientNetB0, all trained to classify ImageNet.

SOFT: Softmax-free Transformer with Linear Complexity

- Jiachen Lu, Jinghan Yao, Junge Zhang, Xiatian Zhu, Hang Xu, Weiguo Gao, Chunjing XU, Tao Xiang, Li Zhang
- abstract@[open-review](#): Vision transformers (ViTs) have pushed the state-of-the-art for various visual recognition tasks by patch-wise image tokenization followed by self-attention. However, the employment of self-attention modules results in a quadratic complexity in both computation and memory usage. Various attempts on approximating the self-attention computation with linear complexity have been made in Natural Language Processing. However, an in-depth analysis in this work shows that they are either theoretically flawed or empirically ineffective for visual recognition. We further identify that their limitations are rooted in keeping the softmax self-attention during approximations. Specifically, conventional self-attention is computed by normalizing the scaled dot-product between token feature vectors. Keeping this softmax operation challenges any subsequent linearization efforts. Based on this insight, for the first time, a softmax-free transformer or SOFT is proposed. To remove softmax in self-attention, Gaussian kernel function is used to replace the dot-product similarity without further normalization. This enables a full self-attention matrix to be approximated via a low-rank matrix decomposition. The robustness of the approximation is achieved by calculating its Moore-Penrose inverse using a Newton-Raphson method. Extensive experiments on ImageNet show that our SOFT significantly improves the computational efficiency of existing ViT variants. Crucially, with a linear complexity, much longer token sequences are permitted in SOFT, resulting in superior trade-off between accuracy and complexity.

Task-Adaptive Neural Network Search with Meta-Contrastive Learning

- Wonyong Jeong, Hayeon Lee, Geon Park, Eunyoung Hyung, Jinheon Baek, Sung Ju Hwang
- abstract@[open-review](#): Most conventional Neural Architecture Search (NAS) approaches are limited in that they only generate architectures without searching for the optimal parameters. While some NAS methods handle this issue by utilizing a supernet trained on a large-scale dataset such as ImageNet, they may be suboptimal if the target tasks are highly dissimilar from the dataset the supernet is trained on. To address such limitations, we introduce a novel problem of Neural Network Search (NNS), whose goal is to search for the optimal pretrained network for a novel dataset and constraints (e.g. number of parameters), from a model zoo. Then, we propose a novel framework to tackle the problem, namely Task-Adaptive Neural Network Search (TANS). Given a model-zoo that consists of network pretrained on diverse datasets, we use a novel amortized meta-learning framework to learn a cross-modal latent space with contrastive loss, to maximize the similarity between a dataset and a high-performing network on it, and minimize the similarity between irrelevant dataset-network pairs. We validate the effectiveness and efficiency of our method on ten real-world datasets, against existing NAS/AutoML baselines. The results show that our method instantly retrieves networks that outperform models obtained with the baselines with significantly fewer training steps to reach the target performance, thus minimizing the total cost of obtaining a task-optimal network. Our code and the model-zoo are available at <https://anonymous.4open.science/r/TANS-33D6>

[Neural Flows: Efficient Alternative to Neural ODEs](#)

- Marin Biloš, Johanna Sommer, Syama Sundar Rangapuram, Tim Januschowski, Stephan Günnemann
- abstract@[open-review](#): Neural ordinary differential equations describe how values change in time. This is the reason why they gained importance in modeling sequential data, especially when the observations are made at irregular intervals. In this paper we propose an alternative by directly modeling the solution curves - the flow of an ODE - with a neural network. This immediately eliminates the need for expensive numerical solvers while still maintaining the modeling capability of neural ODEs. We propose several flow architectures suitable for different applications by establishing precise conditions on when a function defines a valid flow. Apart from computational efficiency, we also provide empirical evidence of favorable generalization performance via applications in time series modeling, forecasting, and density estimation.

[Multi-Objective Meta Learning](#)

- Feiyang YE, Baijiong Lin, Zhixiong Yue, Pengxin Guo, Qiao Xiao, Yu Zhang
- abstract@[open-review](#): Meta learning with multiple objectives has been attracted much attention recently since many applications need to consider multiple factors when designing learning models. Existing gradient-based works on meta learning with multiple objectives mainly combine multiple objectives into a single objective in a weighted sum manner. This simple strategy usually works but it requires to tune the weights associated with all the objectives, which could be time consuming. Different from those works, in this paper, we propose a gradient-based Multi-Objective Meta Learning (MOML) framework without manually tuning weights. Specifically, MOML formulates the objective function of meta learning with multiple objectives as a Multi-Objective Bi-Level optimization Problem (MOBLP) where the upper-level subproblem is to solve several possibly conflicting objectives for the meta learner. To solve the MOBLP, we devise the first gradient-based optimization algorithm by alternatively solving the lower-level and upper-level subproblems via the gradient descent method and the gradient-based multi-objective optimization method, respectively. Theoretically, we prove the convergence properties of the proposed gradient-based optimization algorithm. Empirically, we show the effectiveness of the proposed MOML framework in several meta learning problems, including few-shot learning, domain adaptation, multi-task learning, and neural architecture search. The source code of MOML is available at <https://github.com/Baijiong-Lin/MOML>.

[A self-consistent theory of Gaussian Processes captures feature learning effects in finite CNNs](#)

- Gadi Naveh, Zohar Ringel
- abstract@[open-review](#): Deep neural networks (DNNs) in the infinite width/channel limit have received much attention recently, as they provide a clear analytical window to deep learning via mappings to Gaussian Processes (GPs). Despite its theoretical appeal, this viewpoint lacks a crucial ingredient of deep learning in finite DNNs, laying at the heart of their success --- \textit{feature learning}. Here we consider DNNs trained with noisy gradient descent on a large training set and derive a self-consistent Gaussian Process theory accounting for \textit{strong} finite-DNN and feature learning effects. Applying this to a toy model of a two-layer linear convolutional neural network (CNN) shows good agreement with experiments. We further identify, both analytically and numerically, a sharp transition between a feature learning regime and a lazy learning regime in this model. Strong finite-DNN effects are also derived for a non-linear two-layer fully connected network. We have numerical evidence demonstrating that the assumptions required for our theory hold true in more realistic settings (Myrtle5 CNN trained on CIFAR-10). Our self-consistent theory provides a rich and versatile analytical framework for studying strong finite-DNN effects, most notably - feature learning.

[Mini-Batch Consistent Slot Set Encoder for Scalable Set Encoding](#)

- Andreis Bruno, Jeffrey Willette, Juho Lee, Sung Ju Hwang
- abstract@[open-review](#): Most existing set encoding algorithms operate under the implicit assumption that all the set elements are accessible, and that there are ample computational and memory resources to load the set into memory during training and inference. However, both assumptions fail when the set is excessively large such that it is impossible to load all set elements into memory, or when data arrives in a stream. To tackle such practical challenges in large-scale set encoding, the general set-function constraints of permutation invariance and equivariance are not sufficient. We introduce a new property termed Mini-Batch Consistency (MBC) that is required for large scale mini-batch set encoding. Additionally, we present a scalable and efficient attention-based set encoding mechanism that is amenable to mini-batch processing of sets, and capable of updating set representations as data arrives. The proposed method adheres to the required symmetries of invariance and equivariance as well as maintaining MBC for any partition of the input set. We perform extensive experiments and show that our method is computationally efficient and results in rich set encoding representations for set-structured data.

[Efficient and Local Parallel Random Walks](#)

- Michael Kapralov, Silvio Lattanzi, Navid Nouri, Jakab Tardos
- abstract@[open-review](#): Random walks are a fundamental primitive used in many machine learning algorithms with several applications in clustering and semi-supervised learning. Despite their relevance, the first efficient parallel algorithm to compute random walks has been introduced very recently (Łącki et al.). Unfortunately their method has a fundamental shortcoming: their algorithm is non-local in that it heavily relies on computing random walks out of all nodes in the input graph, even though in many practical applications one is interested in computing random walks only from a small subset of nodes in the graph. In this paper, we present a new algorithm that overcomes this limitation by building random walks efficiently and locally at the same time. We show that our technique is both memory and round efficient, and in particular yields an efficient parallel local clustering algorithm. Finally, we complement our theoretical analysis with experimental results showing that our algorithm is significantly more scalable than previous approaches.

[Amortized Variational Inference for Simple Hierarchical Models](#)

- Abhinav Agrawal, Justin Domke
- abstract@[open-review](#): It is difficult to use subsampling with variational inference in hierarchical models since the number of local latent variables scales with the dataset. Thus, inference in hierarchical models remains a challenge at a large scale. It is helpful to use a variational family with a structure matching the posterior, but optimization is still slow due to the huge number of local distributions. Instead, this paper suggests an amortized approach where shared parameters simultaneously represent all local distributions. This approach is similarly accurate as using a given joint distribution (e.g., a full-

rank Gaussian) but is feasible on datasets that are several orders of magnitude larger. It is also dramatically faster than using a structured variational distribution.

[Online Matching in Sparse Random Graphs: Non-Asymptotic Performances of Greedy Algorithm](#)

- Nathan Noiry, Vianney Perchet, Flore Sentenac
- abstract@[open-review](#): Motivated by sequential budgeted allocation problems, we investigate online matching problems where connections between vertices are not i.i.d., but they have fixed degree distributions -- the so-called configuration model. We estimate the competitive ratio of the simplest algorithm, GREEDY, by approximating some relevant stochastic discrete processes by their continuous counterparts, that are solutions of an explicit system of partial differential equations. This technique gives precise bounds on the estimation errors, with arbitrarily high probability as the problem size increases. In particular, it allows the formal comparison between different configuration models. We also prove that, quite surprisingly, GREEDY can have better performance guarantees than RANKING, another celebrated algorithm for online matching that usually outperforms the former.

[End-to-end reconstruction meets data-driven regularization for inverse problems](#)

- Subhadip Mukherjee, Marcello Carioni, Ozan Öktem, Carola-Bibiane Schönlieb
- abstract@[open-review](#): We propose a new approach for learning end-to-end reconstruction operators based on unpaired training data for ill-posed inverse problems. The proposed method combines the classical variational framework with iterative unrolling and essentially seeks to minimize a weighted combination of the expected distortion in the measurement space and the Wasserstein-1 distance between the distributions of the reconstruction and the ground-truth. More specifically, the regularizer in the variational setting is parametrized by a deep neural network and learned simultaneously with the unrolled reconstruction operator. The variational problem is then initialized with the output of the reconstruction network and solved iteratively till convergence. Notably, it takes significantly fewer iterations to converge as compared to variational methods, thanks to the excellent initialization obtained via the unrolled operator. The resulting approach combines the computational efficiency of end-to-end unrolled reconstruction with the well-posedness and noise-stability guarantees of the variational setting. Moreover, we demonstrate with the example of image reconstruction in X-ray computed tomography (CT) that our approach outperforms state-of-the-art unsupervised methods and that it outperforms or is at least on par with state-of-the-art supervised data-driven reconstruction approaches.

[An online passive-aggressive algorithm for difference-of-squares classification](#)

- Lawrence Saul
- abstract@[open-review](#): We investigate a low-rank model of quadratic classification inspired by previous work on factorization machines, polynomial networks, and capsule-based architectures for visual object recognition. The model is parameterized by a pair of affine transformations, and it classifies examples by comparing the magnitudes of vectors that these transformations produce. The model is also over-parameterized in the sense that different pairs of affine transformations can describe classifiers with the same decision boundary and confidence scores. We show that such pairs arise from discrete and continuous symmetries of the model's parameter space: in particular, the latter define symmetry groups of rotations and Lorentz transformations, and we use these group structures to devise appropriately invariant procedures for model alignment and averaging. We also leverage the form of the model's decision boundary to derive simple margin-based updates for online learning. Here we explore a strategy of passive-aggressive learning: for each example, we compute the minimum change in parameters that is required to predict its correct label with high confidence. We derive these updates by solving a quadratically constrained quadratic program (QCQP); interestingly, this QCQP is nonconvex but tractable, and it can be solved efficiently by elementary methods. We highlight the conceptual and practical contributions of this approach. Conceptually, we show that it extends the paradigm of passive-aggressive learning to a larger family of nonlinear models for classification. Practically, we show that these models perform well on large-scale problems in online learning.

[Finite-Sample Analysis of Off-Policy TD-Learning via Generalized Bellman Operators](#)

- Zaiwei Chen, Siva Theja Maguluri, Sanjay Shakkottai, Karthikeyan Shanmugam
- abstract@[open-review](#): In TD-learning, off-policy sampling is known to be more practical than on-policy sampling, and by decoupling learning from data collection, it enables data reuse. It is known that policy evaluation has the interpretation of solving a generalized Bellman equation. In this paper, we derive finite-sample bounds for any general off-policy TD-like stochastic approximation algorithm that solves for the fixed-point of this generalized Bellman operator. Our key step is to show that the generalized Bellman operator is simultaneously a contraction mapping with respect to a weighted $\|\cdot\|_p$ -norm for each $p \in [1, \infty)$, with a common contraction factor. Off-policy TD-learning is known to suffer from high variance due to the product of importance sampling ratios. A number of algorithms (e.g. $Q^\pi(\lambda)$, Tree-Backup (λ) , Retrace (λ) , and Q -trace) have been proposed in the literature to address this issue. Our results immediately imply finite-sample bounds of these algorithms. In particular, we provide first-known finite-sample guarantees for $Q^\pi(\lambda)$, Tree-Backup (λ) , and Retrace (λ) , and improve the best known bounds of Q -trace in [chen2021finite](#). Moreover, we show the bias-variance trade-offs in each of these algorithms.

[A Bi-Level Framework for Learning to Solve Combinatorial Optimization on Graphs](#)

- Runzhong Wang, Zhigang Hua, Gan Liu, Jiayi Zhang, Junchi Yan, Feng Qi, Shuang Yang, Jun Zhou, Xiaokang Yang
- abstract@[open-review](#): Combinatorial Optimization (CO) has been a long-standing challenging research topic featured by its NP-hard nature. Traditionally such problems are approximately solved with heuristic algorithms which are usually fast but may sacrifice the solution quality. Currently, machine learning for combinatorial optimization (MLCO) has become a trending research topic, but most existing MLCO methods treat CO as a single-level optimization by directly learning the end-to-end solutions, which are hard to scale up and mostly limited by the capacity of ML models given the high complexity of CO. In this paper, we propose a hybrid approach to combine the best of the two worlds, in which a bi-level framework is developed with an upper-level learning method to optimize the graph (e.g. add, delete or modify edges in a graph), fused with a lower-level heuristic algorithm solving on the optimized graph. Such a bi-level approach simplifies the learning on the original hard CO and can effectively mitigate the demand for model capacity. The experiments and results on several popular CO problems like Directed Acyclic Graph scheduling, Graph Edit Distance and Hamiltonian Cycle Problem show its effectiveness over manually designed heuristics and single-level learning methods.

[Improved Learning Rates of a Functional Lasso-type SVM with Sparse Multi-Kernel Representation](#)

- shaogao lv, Junhui Wang, Jiankun Liu, Yong Liu
- abstract@[open-review](#): In this paper, we provide theoretical results of estimation bounds and excess risk upper bounds for support vector machine (SVM) with sparse multi-kernel representation. These convergence rates for multi-kernel SVM are established by analyzing a Lasso-type regularized learning scheme within composite multi-kernel spaces. It is shown that the oracle rates of convergence of classifiers depend on the complexity of multi-kernels, the sparsity, a Bernstein condition and the sample size, which significantly improves on previous results even for the additive or linear cases. In summary, this paper not only provides unified theoretical results for multi-kernel SVMs, but also enriches the literature on high-dimensional nonparametric classification.

[When does Contrastive Learning Preserve Adversarial Robustness from Pretraining to Finetuning?](#)

- Lijie Fan, Sijia Liu, Pin-Yu Chen, Gaoyuan Zhang, Chuang Gan
- abstract@[open-review](#): Contrastive learning (CL) can learn generalizable feature representations and achieve state-of-the-art performance of downstream tasks by finetuning a linear classifier on top of it. However, as adversarial robustness becomes vital in image classification, it remains unclear whether or not CL is able to preserve robustness to downstream tasks. The main challenge is that in the self-supervised pretraining + supervised finetuning paradigm, adversarial robustness is easily forgotten due to a learning task mismatch from pretraining to finetuning. We call such challenge 'cross-task robustness transferability'. To address the above problem, in this paper we revisit and advance CL principles through the lens of robustness enhancement. We show that (1) the design of contrastive views matters: High-frequency components of images are beneficial to improving model robustness; (2) Augmenting CL with pseudo-supervision stimulus (e.g., resorting to feature clustering) helps preserve robustness without forgetting. Equipped with our new designs, we propose AdvCL, a novel adversarial contrastive pretraining framework. We show that AdvCL is able to enhance cross-task robustness transferability without loss of model accuracy and finetuning efficiency. With a thorough experimental study, we demonstrate that AdvCL outperforms the state-of-the-art self-supervised robust learning methods across multiple datasets (CIFAR-10, CIFAR-100, and STL-10) and finetuning schemes (linear evaluation and full model finetuning).

[Learning Transferable Features for Point Cloud Detection via 3D Contrastive Co-training](#)

- Zeng Yihan, Chunwei Wang, Yunbo Wang, Hang Xu, Chaoqiang Ye, Zhen Yang, Chao Ma
- abstract@[open-review](#): Most existing point cloud detection models require large-scale, densely annotated datasets. They typically underperform in domain adaptation settings, due to geometry shifts caused by different physical environments or LiDAR sensor configurations. Therefore, it is challenging but valuable to learn transferable features between a labeled source domain and a novel target domain, without any access to target labels. To tackle this problem, we introduce the framework of 3D Contrastive Co-training (3D-CoCo) with two technical contributions. First, 3D-CoCo is inspired by our observation that the bird-eye-view (BEV) features are more transferable than low-level geometry features. We thus propose a new co-training architecture that includes separate 3D encoders with domain-specific parameters, as well as a BEV transformation module for learning domain-invariant features. Second, 3D-CoCo extends the approach of contrastive instance alignment to point cloud detection, whose performance was largely hindered by the mismatch between the fictitious distribution of BEV features, induced by pseudo-labels, and the true distribution. The mismatch is greatly reduced by 3D-CoCo with transformed point clouds, which are carefully designed by considering specific geometry priors. We construct new domain adaptation benchmarks using three large-scale 3D datasets. Experimental results show that our proposed 3D-CoCo effectively closes the domain gap and outperforms the state-of-the-art methods by large margins.

[SILG: The Multi-domain Symbolic Interactive Language Grounding Benchmark](#)

- Victor Zhong, Austin W. Hanjie, Sida Wang, Karthik Narasimhan, Luke Zettlemoyer
- abstract@[open-review](#): Existing work in language grounding typically study single environments. How do we build unified models that apply across multiple environments? We propose the multi-environment Symbolic Interactive Language Grounding benchmark (SILG), which unifies a collection of diverse grounded language learning environments under a common interface. SILG consists of grid-world environments that require generalization to new dynamics, entities, and partially observed worlds (RTFM, Messenger, NetHack), as well as symbolic counterparts of visual worlds that require interpreting rich natural language with respect to complex scenes (ALFWORLD, Touchdown). Together, these environments provide diverse grounding challenges in richness of observation space, action space, language specification, and plan complexity. In addition, we propose the first shared model architecture for RL on these environments, and evaluate recent advances such as egocentric local convolution, recurrent state-tracking, entity-centric attention, and pretrained LM using SILG. Our shared architecture achieves comparable performance to environment-specific architectures. Moreover, we find that many recent modelling advances do not result in significant gains on environments other than the one they were designed for. This highlights the need for a multi-environment benchmark. Finally, the best models significantly underperform humans on SILG, which suggests ample room for future work. We hope SILG enables the community to quickly identify new methodologies for language grounding that generalize to a diverse set of environments and their associated challenges.

[A Surrogate Objective Framework for Prediction+Programming with Soft Constraints](#)

- Kai Yan, Jie Yan, Chuan Luo, Liting Chen, Qingwei Lin, Dongmei Zhang
- abstract@[open-review](#): Prediction+optimization is a common real-world paradigm where we have to predict problem parameters before solving the optimization problem. However, the criteria by which the prediction model is trained are often inconsistent with the goal of the downstream optimization problem. Recently, decision-focused prediction approaches, such as SPO+ and direct optimization, have been proposed to fill this gap. However, they cannot directly handle the soft constraints with the max operator required in many real-world objectives. This paper proposes a novel analytically differentiable surrogate objective framework for real-world linear and semi-definite negative quadratic programming problems with soft linear and non-negative hard constraints. This framework gives the theoretical bounds on constraints' multipliers, and derives the closed-form solution with respect to predictive parameters and thus gradients for any variable in the problem. We evaluate our method in three applications extended with soft constraints: synthetic linear programming, portfolio optimization, and resource provisioning, demonstrating that our method outperforms traditional two-staged methods and other decision-focused approaches

[Learning to Predict Trustworthiness with Steep Slope Loss](#)

- Yan Luo, Yongkang Wong, Mohan S. Kankanhalli, Qi Zhao
- abstract@[open-review](#): Understanding the trustworthiness of a prediction yielded by a classifier is critical for the safe and effective use of AI models. Prior efforts have been proven to be reliable on small-scale datasets. In this work, we study the problem of predicting trustworthiness on real-world large-scale datasets, where the task is more challenging due to high-dimensional features, diverse visual concepts, and a large number of samples. In such a setting, we observe that the trustworthiness predictors trained with prior-art loss functions, i.e., the cross entropy loss, focal loss, and true class probability confidence loss, are prone to view both correct predictions and incorrect predictions to be trustworthy. The reasons are two-fold. Firstly, correct predictions are generally dominant over incorrect predictions. Secondly, due to the data complexity, it is challenging to differentiate the incorrect predictions from the correct ones on real-world large-scale datasets. To improve the generalizability of trustworthiness predictors, we propose a novel steep slope loss to separate the features w.r.t. correct predictions from the ones w.r.t. incorrect predictions by two slide-like curves that oppose each other. The proposed loss is evaluated with two representative deep learning models, i.e., Vision Transformer and ResNet, as trustworthiness predictors. We conduct comprehensive experiments and analyses on ImageNet, which show that the proposed loss effectively improves the generalizability of trustworthiness predictors. The code and pre-trained trustworthiness predictors for reproducibility are available at \url{https://github.com/luoyan407/predict_trustworthiness}.

[On the Periodic Behavior of Neural Network Training with Batch Normalization and Weight Decay](#)

- Ekaterina Lobacheva, Maxim Kodryan, Nadezhda Chirkova, Andrey Malinin, Dmitry P. Vetrov
- abstract@[open-review](#): Training neural networks with batch normalization and weight decay has become a common practice in recent years. In this work, we show that their combined use may result in a surprising periodic behavior of optimization dynamics: the training process regularly exhibits destabilizations that, however, do not lead to complete divergence but cause a new period of training. We rigorously investigate the mechanism underlying the discovered periodic behavior from both empirical and theoretical points of view and analyze the conditions in which it occurs in practice. We also demonstrate that periodic behavior can be regarded as a generalization of two previously opposing perspectives on training with batch normalization and weight decay, namely the equilibrium presumption and the instability presumption.

NeRV: Neural Representations for Videos

- Hao Chen, Bo He, Hanyu Wang, Yixuan Ren, Ser Nam Lim, Abhinav Shrivastava
- abstract@[open-review](#): We propose a novel neural representation for videos (NeRV) which encodes videos in neural networks. Unlike conventional representations that treat videos as frame sequences, we represent videos as neural networks taking frame index as input. Given a frame index, NeRV outputs the corresponding RGB image. Video encoding in NeRV is simply fitting a neural network to video frames and decoding process is a simple feedforward operation. As an image-wise implicit representation, NeRV output the whole image and shows great efficiency compared to pixel-wise implicit representation, improving the encoding speed by \$textbf{25}\times\$ to \$textbf{70}\times\$, the decoding speed by \$textbf{38}\times\$ to \$textbf{132}\times\$, while achieving better video quality. With such a representation, we can treat videos as neural networks, simplifying several video-related tasks. For example, conventional video compression methods are restricted by a long and complex pipeline, specifically designed for the task. In contrast, with NeRV, we can use any neural network compression method as a proxy for video compression, and achieve comparable performance to traditional frame-based video compression approaches (H.264, HEVC \etc). Besides compression, we demonstrate the generalization of NeRV for video denoising. The source code and pre-trained model can be found at <https://github.com/haochen-rye/NeRV.git>.

Surrogate Regret Bounds for Polyhedral Losses

- Rafael Frongillo, Bo Waggoner
- abstract@[open-review](#): Surrogate risk minimization is an ubiquitous paradigm in supervised machine learning, wherein a target problem is solved by minimizing a surrogate loss on a dataset. Surrogate regret bounds, also called excess risk bounds, are a common tool to prove generalization rates for surrogate risk minimization. While surrogate regret bounds have been developed for certain classes of loss functions, such as proper losses, general results are relatively sparse. We provide two general results. The first gives a linear surrogate regret bound for any polyhedral (piecewise-linear and convex) surrogate, meaning that surrogate generalization rates translate directly to target rates. The second shows that for sufficiently non-polyhedral surrogates, the regret bound is a square root, meaning fast surrogate generalization rates translate to slow rates for the target. Together, these results suggest polyhedral surrogates are optimal in many cases.

Last iterate convergence of SGD for Least-Squares in the Interpolation regime.

- Aditya Vardhan Varre, Loucas Pillaud-Vivien, Nicolas Flammarion
- abstract@[open-review](#): Motivated by the recent successes of neural networks that have the ability to fit the data perfectly \emph{and} generalize well, we study the noiseless model in the fundamental least-squares setup. We assume that an optimum predictor perfectly fits the inputs and outputs \$\langle \theta^*, \phi(X) \rangle = Y\$, where \$\phi(X)\$ stands for a possibly infinite dimensional non-linear feature map. To solve this problem, we consider the estimator given by the last iterate of stochastic gradient descent (SGD) with constant step-size. In this context, our contribution is two fold: (i) \emph{from a (stochastic) optimization perspective}, we exhibit an archetypal problem where we can show explicitly the convergence of SGD final iterate for a non-strongly convex problem with constant step-size whereas usual results use some form of average and (ii) \emph{from a statistical perspective}, we give explicit non-asymptotic convergence rates in the over-parameterized setting and leverage a \emph{fine-grained} parameterization of the problem to exhibit polynomial rates that can be faster than \$O(1/T)\$. The link with reproducing kernel Hilbert spaces is established.

Generative vs. Discriminative: Rethinking The Meta-Continual Learning

- Mohammadamin Banayeeanzade, Rasoul Mirzaiezadeh, Hosein Hasani, Mahdieh Soleymani
- abstract@[open-review](#): Deep neural networks have achieved human-level capabilities in various learning tasks. However, they generally lose performance in more realistic scenarios like learning in a continual manner. In contrast, humans can incorporate their prior knowledge to learn new concepts efficiently without forgetting older ones. In this work, we leverage meta-learning to encourage the model to learn how to learn continually. Inspired by human concept learning, we develop a generative classifier that efficiently uses data-driven experience to learn new concepts even from few samples while being immune to forgetting. Along with cognitive and theoretical insights, extensive experiments on standard benchmarks demonstrate the effectiveness of the proposed method. The ability to remember all previous concepts, with negligible computational and structural overheads, suggests that generative models provide a natural way for alleviating catastrophic forgetting, which is a major drawback of discriminative models.

Model, sample, and epoch-wise descents: exact solution of gradient flow in the random feature model

- Antoine Bodin, Nicolas Macris
- abstract@[open-review](#): Recent evidence has shown the existence of a so-called double-descent and even triple-descent behavior for the generalization error of deep-learning models. This important phenomenon commonly appears in implemented neural network architectures, and also seems to emerge in epoch-wise curves during the training process. A recent line of research has highlighted that random matrix tools can be used to obtain precise analytical asymptotics of the generalization (and training) errors of the random feature model. In this contribution, we analyze the whole temporal behavior of the generalization and training errors under gradient flow for the random feature model. We show that in the asymptotic limit of large system size the full time-evolution path of both errors can be calculated analytically. This allows us to observe how the double and triple descents develop over time, if and when early stopping is an option, and also observe time-wise descent structures. Our techniques are based on Cauchy complex integral representations of the errors together with recent random matrix methods based on linear pencils.

Rethinking Graph Transformers with Spectral Attention

- Devin Kreuzer, Dominique Beaini, Will Hamilton, Vincent Létourneau, Prudencio Tossou
- abstract@[open-review](#): In recent years, the Transformer architecture has proven to be very successful in sequence processing, but its application to other data structures, such as graphs, has remained limited due to the difficulty of properly defining positions. Here, we present the \textit{Spectral Attention Network} (SAN), which uses a learned positional encoding (LPE) that can take advantage of the full Laplacian spectrum to learn the position of each node in a given graph. This LPE is then added to the node features of the graph and passed to a fully-connected Transformer. By leveraging the full spectrum of the Laplacian, our model is theoretically powerful in distinguishing graphs, and can better detect similar sub-structures from their resonance. Further, by fully connecting the graph, the Transformer does not suffer from over-squashing, an information bottleneck of most GNNs, and enables better modeling of physical phenomena such as heat transfer and electric interaction. When tested empirically on a set of 4 standard datasets, our model performs on par or better than state-of-the-art GNNs, and outperforms any attention-based model by a wide margin, becoming the first fully-connected architecture to perform well on graph benchmarks.

Perceptual Score: What Data Modalities Does Your Model Perceive?

- Itai Gat, Idan Schwartz, Alex Schwing
- abstract@[open-review](#): Machine learning advances in the last decade have relied significantly on large-scale datasets that continue to grow in size. Increasingly, those datasets also contain different data modalities. However, large multi-modal datasets are hard to annotate, and annotations may contain biases that we are often unaware of. Deep-net-based classifiers, in turn, are prone to exploit those biases and to find shortcuts. To study and quantify this concern, we introduce the perceptual score, a metric that assesses the degree to which a model relies on the different subsets of the input features, i.e., modalities. Using the perceptual score, we find a surprisingly consistent trend across four popular datasets: recent, more accurate state-of-the-art multi-

modal models for visual question-answering or visual dialog tend to perceive the visual data less than their predecessors. This is concerning as answers are hence increasingly inferred from textual cues only. Using the perceptual score also helps to analyze model biases by decomposing the score into data subset contributions. We hope to spur a discussion on the perceptiveness of multi-modal models and also hope to encourage the community working on multi-modal classifiers to start quantifying perceptiveness via the proposed perceptual score.

[PiRank: Scalable Learning To Rank via Differentiable Sorting](#)

- Robin Swezey, Aditya Grover, Bruno Charron, Stefano Ermon
- abstract@[open-review](#): A key challenge with machine learning approaches for ranking is the gap between the performance metrics of interest and the surrogate loss functions that can be optimized with gradient-based methods. This gap arises because ranking metrics typically involve a sorting operation which is not differentiable w.r.t. the model parameters. Prior works have proposed surrogates that are loosely related to ranking metrics or simple smoothed versions thereof, and often fail to scale to real-world applications. We propose PiRank, a new class of differentiable surrogates for ranking, which employ a continuous, temperature-controlled relaxation to the sorting operator based on NeuralSort [1]. We show that PiRank exactly recovers the desired metrics in the limit of zero temperature and further propose a divide-and-conquer extension that scales favorably to large list sizes, both in theory and practice. Empirically, we demonstrate the role of larger list sizes during training and show that PiRank significantly improves over comparable approaches on publicly available Internet-scale learning-to-rank benchmarks.

[Deceive D: Adaptive Pseudo Augmentation for GAN Training with Limited Data](#)

- Liming Jiang, Bo Dai, Wayne Wu, Chen Change Loy
- abstract@[open-review](#): Generative adversarial networks (GANs) typically require ample data for training in order to synthesize high-fidelity images. Recent studies have shown that training GANs with limited data remains formidable due to discriminator overfitting, the underlying cause that impedes the generator's convergence. This paper introduces a novel strategy called Adaptive Pseudo Augmentation (APA) to encourage healthy competition between the generator and the discriminator. As an alternative method to existing approaches that rely on standard data augmentations or model regularization, APA alleviates overfitting by employing the generator itself to augment the real data distribution with generated images, which deceives the discriminator adaptively. Extensive experiments demonstrate the effectiveness of APA in improving synthesis quality in the low-data regime. We provide a theoretical analysis to examine the convergence and rationality of our new training strategy. APA is simple and effective. It can be added seamlessly to powerful contemporary GANs, such as StyleGAN2, with negligible computational cost. Code: <https://github.com/EndlessSora/DeceiveD>.

[CoFrNets: Interpretable Neural Architecture Inspired by Continued Fractions](#)

- Isha Puri, Amit Dhurandhar, Tejaswini Pedapati, Karthikeyan Shanmugam, Dennis Wei, Kush R. Varshney
- abstract@[open-review](#): In recent years there has been a considerable amount of research on local post hoc explanations for neural networks. However, work on building interpretable neural architectures has been relatively sparse. In this paper, we present a novel neural architecture, CoFrNet, inspired by the form of continued fractions which are known to have many attractive properties in number theory, such as fast convergence of approximations to real numbers. We show that CoFrNets can be efficiently trained as well as interpreted leveraging their particular functional form. Moreover, we prove that such architectures are universal approximators based on a proof strategy that is different than the typical strategy used to prove universal approximation results for neural networks based on infinite width (or depth), which is likely to be of independent interest. We experiment on nonlinear synthetic functions and are able to accurately model as well as estimate feature attributions and even higher order terms in some cases, which is a testament to the representational power as well as interpretability of such architectures. To further showcase the power of CoFrNets, we experiment on seven real datasets spanning tabular, text and image modalities, and show that they are either comparable or significantly better than other interpretable models and multilayer perceptrons, sometimes approaching the accuracies of state-of-the-art models.

[Iterative Teaching by Label Synthesis](#)

- Weiyang Liu, Zhen Liu, Hanchen Wang, Liam Paull, Bernhard Schölkopf, Adrian Weller
- abstract@[open-review](#): In this paper, we consider the problem of iterative machine teaching, where a teacher provides examples sequentially based on the current iterative learner. In contrast to previous methods that have to scan over the entire pool and select teaching examples from it in each iteration, we propose a label synthesis teaching framework where the teacher randomly selects input teaching examples (e.g., images) and then synthesizes suitable outputs (e.g., labels) for them. We show that this framework can avoid costly example selection while still provably achieving exponential teachability. We propose multiple novel teaching algorithms in this framework. Finally, we empirically demonstrate the value of our framework.

[Variational Diffusion Models](#)

- Diederik Kingma, Tim Salimans, Ben Poole, Jonathan Ho
- abstract@[open-review](#): Diffusion-based generative models have demonstrated a capacity for perceptually impressive synthesis, but can they also be great likelihood-based models? We answer this in the affirmative, and introduce a family of diffusion-based generative models that obtain state-of-the-art likelihoods on standard image density estimation benchmarks. Unlike other diffusion-based models, our method allows for efficient optimization of the noise schedule jointly with the rest of the model. We show that the variational lower bound (VLB) simplifies to a remarkably short expression in terms of the signal-to-noise ratio of the diffused data, thereby improving our theoretical understanding of this model class. Using this insight, we prove an equivalence between several models proposed in the literature. In addition, we show that the continuous-time VLB is invariant to the noise schedule, except for the signal-to-noise ratio at its endpoints. This enables us to learn a noise schedule that minimizes the variance of the resulting VLB estimator, leading to faster optimization. Combining these advances with architectural improvements, we obtain state-of-the-art likelihoods on image density estimation benchmarks, outperforming autoregressive models that have dominated these benchmarks for many years, with often significantly faster optimization. In addition, we show how to use the model as part of a bits-back compression scheme, and demonstrate lossless compression rates close to the theoretical optimum.

[FastCorrect: Fast Error Correction with Edit Alignment for Automatic Speech Recognition](#)

- Yichong Leng, Xu Tan, Linchen Zhu, Jin Xu, Renqian Luo, Linquan Liu, Tao Qin, Xiangyang Li, Edward Lin, Tie-Yan Liu
- abstract@[open-review](#): Error correction techniques have been used to refine the output sentences from automatic speech recognition (ASR) models and achieve a lower word error rate (WER) than original ASR outputs. Previous works usually use a sequence-to-sequence model to correct an ASR output sentence autoregressively, which causes large latency and cannot be deployed in online ASR services. A straightforward solution to reduce latency, inspired by non-autoregressive (NAR) neural machine translation, is to use an NAR sequence generation model for ASR error correction, which, however, comes at the cost of significantly increased ASR error rate. In this paper, observing distinctive error patterns and correction operations (i.e., insertion, deletion, and substitution) in ASR, we propose FastCorrect, a novel NAR error correction model based on edit alignment. In training, FastCorrect aligns each source token from an ASR output sentence to the target tokens from the corresponding ground-truth sentence based on the edit distance between the source and target sentences, and extracts the number of target tokens corresponding to each source token during edition/correction, which is then used to train a length predictor and to adjust the source tokens to match the length of the target sentence for parallel generation. In inference, the token number predicted by the length predictor is used to adjust the source tokens for target sequence generation. Experiments on the public AISHELL-1 dataset and an internal industrial-scale ASR dataset show the effectiveness of FastCorrect for ASR error correction: 1) it speeds up the inference by 6-9 times and

maintains the accuracy (8-14% WER reduction) compared with the autoregressive correction model; and 2) it outperforms the popular NAR models adopted in neural machine translation and text edition by a large margin.

[Integrated Latent Heterogeneity and Invariance Learning in Kernel Space](#)

- Jiashuo Liu, Zheyuan Hu, Peng Cui, Bo Li, Zheyuan Shen
- abstract@[open-review](#): The ability to generalize under distributional shifts is essential to reliable machine learning, while models optimized with empirical risk minimization usually fail on non-\$i.i.d\$ testing data. Recently, invariant learning methods for out-of-distribution (OOD) generalization propose to find causally invariant relationships with multi-environments. However, modern datasets are frequently multi-sourced without explicit source labels, rendering many invariant learning methods inapplicable. In this paper, we propose Kernelized Heterogeneous Risk Minimization (KerHRM) algorithm, which achieves both the latent heterogeneity exploration and invariant learning in kernel space, and then gives feedback to the original neural network by appointing invariant gradient direction. We theoretically justify our algorithm and empirically validate the effectiveness of our algorithm with extensive experiments.

[Hierarchical Reinforcement Learning with Timed Subgoals](#)

- Nico Gürler, Dieter Büchler, Georg Martius
- abstract@[open-review](#): Hierarchical reinforcement learning (HRL) holds great potential for sample-efficient learning on challenging long-horizon tasks. In particular, letting a higher level assign subgoals to a lower level has been shown to enable fast learning on difficult problems. However, such subgoal-based methods have been designed with static reinforcement learning environments in mind and consequently struggle with dynamic elements beyond the immediate control of the agent even though they are ubiquitous in real-world problems. In this paper, we introduce Hierarchical reinforcement learning with Timed Subgoals (HiTS), an HRL algorithm that enables the agent to adapt its timing to a dynamic environment by not only specifying what goal state is to be reached but also when. We discuss how communicating with a lower level in terms of such timed subgoals results in a more stable learning problem for the higher level. Our experiments on a range of standard benchmarks and three new challenging dynamic reinforcement learning environments show that our method is capable of sample-efficient learning where an existing state-of-the-art subgoal-based HRL method fails to learn stable solutions.

[Fair Scheduling for Time-dependent Resources](#)

- Bo Li, Minming Li, Ruilong Zhang
- abstract@[open-review](#): We study a fair resource scheduling problem, where a set of interval jobs are to be allocated to heterogeneous machines controlled by intellectual agents. Each job is associated with release time, deadline, and processing time such that it can be processed if its complete processing period is between its release time and deadline. The machines gain possibly different utilities by processing different jobs, and all jobs assigned to the same machine should be processed without overlap. We consider two widely studied solution concepts, namely, maximin share fairness and envy-freeness. For both criteria, we discuss the extent to which fair allocations exist and present constant approximation algorithms for various settings.

[SNIPS: Solving Noisy Inverse Problems Stochastically](#)

- Bahjat Kawar, Gregory Vaksman, Michael Elad
- abstract@[open-review](#): In this work we introduce a novel stochastic algorithm dubbed SNIPS, which draws samples from the posterior distribution of any linear inverse problem, where the observation is assumed to be contaminated by additive white Gaussian noise. Our solution incorporates ideas from Langevin dynamics and Newton's method, and exploits a pre-trained minimum mean squared error (MMSE) Gaussian denoiser. The proposed approach relies on an intricate derivation of the posterior score function that includes a singular value decomposition (SVD) of the degradation operator, in order to obtain a tractable iterative algorithm for the desired sampling. Due to its stochasticity, the algorithm can produce multiple high perceptual quality samples for the same noisy observation. We demonstrate the abilities of the proposed paradigm for image deblurring, super-resolution, and compressive sensing. We show that the samples produced are sharp, detailed and consistent with the given measurements, and their diversity exposes the inherent uncertainty in the inverse problem being solved.

[Stateful ODE-Nets using Basis Function Expansions](#)

- Alejandro Queiruga, N. Benjamin Erichson, Liam Hodgkinson, Michael W. Mahoney
- abstract@[open-review](#): The recently-introduced class of ordinary differential equation networks (ODE-Nets) establishes a fruitful connection between deep learning and dynamical systems. In this work, we reconsider formulations of the weights as continuous-in-depth functions using linear combinations of basis functions which enables us to leverage parameter transformations such as function projections. In turn, this view allows us to formulate a novel stateful ODE-Block that handles stateful layers. The benefits of this new ODE-Block are twofold: first, it enables incorporating meaningful continuous-in-depth batch normalization layers to achieve state-of-the-art performance; second, it enables compressing the weights through a change of basis, without retraining, while maintaining near state-of-the-art performance and reducing both inference time and memory footprint. Performance is demonstrated by applying our stateful ODE-Block to (a) image classification tasks using convolutional units and (b) sentence-tagging tasks using transformer encoder units.

[Beyond the Signs: Nonparametric Tensor Completion via Sign Series](#)

- Chanwoo Lee, Miaoyan Wang
- abstract@[open-review](#): We consider the problem of tensor estimation from noisy observations with possibly missing entries. A nonparametric approach to tensor completion is developed based on a new model which we coin as sign representable tensors. The model represents the signal tensor of interest using a series of structured sign tensors. Unlike earlier methods, the sign series representation effectively addresses both low- and high-rank signals, while encompassing many existing tensor models---including CP models, Tucker models, single index models, structured tensors with repeating entries---as special cases. We provably reduce the tensor estimation problem to a series of structured classification tasks, and we develop a learning reduction machinery to empower existing low-rank tensor algorithms for more challenging high-rank estimation. Excess risk bounds, estimation errors, and sample complexities are established. We demonstrate the outperformance of our approach over previous methods on two datasets, one on human brain connectivity networks and the other on topic data mining.

[Functional Variational Inference based on Stochastic Process Generators](#)

- Chao Ma, José Miguel Hernández-Lobato
- abstract@[open-review](#): Bayesian inference in the space of functions has been an important topic for Bayesian modeling in the past. In this paper, we propose a new solution to this problem called Functional Variational Inference (FVI). In FVI, we minimize a divergence in function space between the variational distribution and the posterior process. This is done by using as functional variational family a new class of flexible distributions called Stochastic Process Generators (SPGs), which are cleverly designed so that the functional ELBO can be estimated efficiently using analytic solutions and mini-batch sampling. FVI can be applied to stochastic process priors when random function samples from those priors are available. Our experiments show that FVI consistently outperforms weight-space and function space VI methods on several tasks, which validates the effectiveness of our approach.

[TTT++: When Does Self-Supervised Test-Time Training Fail or Thrive?](#)

- Yuejiang Liu, Parth Kothari, Bastien van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, Alexandre Alahi
- abstract@[open-review](#): Test-time training (TTT) through self-supervised learning (SSL) is an emerging paradigm to tackle distributional shifts. Despite encouraging results, it remains unclear when this approach thrives or fails. In this work, we first provide an in-depth look at its limitations and show that TTT can possibly deteriorate, instead of improving, the test-time performance in the presence of severe distribution shifts. To address this issue, we introduce a test-time feature alignment strategy utilizing offline feature summarization and online moment matching, which regularizes adaptation without revisiting training data. We further scale this strategy in the online setting through batch-queue decoupling to enable robust moment estimates even with limited batch size. Given aligned feature distributions, we then shed light on the strong potential of TTT by theoretically analyzing its performance post adaptation. This analysis motivates our use of more informative self-supervision in the form of contrastive learning for visual recognition problems. We empirically demonstrate that our modified version of test-time training, termed TTT++, outperforms state-of-the-art methods by significant margins on several benchmarks. Our result indicates that storing and exploiting extra information, in addition to model parameters, can be a promising direction towards robust test-time adaptation.

[Double Machine Learning Density Estimation for Local Treatment Effects with Instruments](#)

- Yonghan Jung, Jin Tian, Elias Bareinboim
- abstract@[open-review](#): Local treatment effects are a common quantity found throughout the empirical sciences that measure the treatment effect among those who comply with what they are assigned. Most of the literature is focused on estimating the average of such quantity, which is called the ``local average treatment effect (LATE)'' [Imbens and Angrist, 1994]). In this work, we study how to estimate the density of the local treatment effect, which is naturally more informative than its average. Specifically, we develop two families of methods for this task, namely, kernel-smoothing and model-based approaches. The kernel-smoothing-based approach estimates the density through some smooth kernel functions. The model-based approach estimates the density by projecting it onto a finite-dimensional density class. For both approaches, we derive the corresponding double/debiased machine learning-based estimators [Chernozhukov et al., 2018]. We further study the asymptotic convergence rates of the estimators and show that they are robust to the biases in nuisance function estimation. The use of the proposed methods is illustrated through both synthetic and a real dataset called 401(k).

[Dirichlet Energy Constrained Learning for Deep Graph Neural Networks](#)

- Kaixiong Zhou, Xiao Huang, Daochen Zha, Rui Chen, Li Li, Soo-Hyun Choi, Xia Hu
- abstract@[open-review](#): Graph neural networks (GNNs) integrate deep architectures and topological structure modeling in an effective way. However, the performance of existing GNNs would decrease significantly when they stack many layers, because of the over-smoothing issue. Node embeddings tend to converge to similar vectors when GNNs keep recursively aggregating the representations of neighbors. To enable deep GNNs, several methods have been explored recently. But they are developed from either techniques in convolutional neural networks or heuristic strategies. There is no generalizable and theoretical principle to guide the design of deep GNNs. To this end, we analyze the bottleneck of deep GNNs by leveraging the Dirichlet energy of node embeddings, and propose a generalizable principle to guide the training of deep GNNs. Based on it, a novel deep GNN framework -- Energetic Graph Neural Networks (EGNN) is designed. It could provide lower and upper constraints in terms of Dirichlet energy at each layer to avoid over-smoothing. Experimental results demonstrate that EGNN achieves state-of-the-art performance by using deep layers.

[Accelerating Robotic Reinforcement Learning via Parameterized Action Primitives](#)

- Murtaza Dalal, Deepak Pathak, Russ R. Salakhutdinov
- abstract@[open-review](#): Despite the potential of reinforcement learning (RL) for building general-purpose robotic systems, training RL agents to solve robotics tasks still remains challenging due to the difficulty of exploration in purely continuous action spaces. Addressing this problem is an active area of research with the majority of focus on improving RL methods via better optimization or more efficient exploration. An alternate but important component to consider improving is the interface of the RL algorithm with the robot. In this work, we manually specify a library of robot action primitives (RAPS), parameterized with arguments that are learned by an RL policy. These parameterized primitives are expressive, simple to implement, enable efficient exploration and can be transferred across robots, tasks and environments. We perform a thorough empirical study across challenging tasks in three distinct domains with image input and a sparse terminal reward. We find that our simple change to the action interface substantially improves both the learning efficiency and task performance irrespective of the underlying RL algorithm, significantly outperforming prior methods which learn skills from offline expert data.

[Boosted CVaR Classification](#)

- Runtian Zhai, Chen Dan, Arun Suggala, J. Zico Kolter, Pradeep Ravikumar
- abstract@[open-review](#): Many modern machine learning tasks require models with high tail performance, i.e. high performance over the worst-off samples in the dataset. This problem has been widely studied in fields such as algorithmic fairness, class imbalance, and risk-sensitive decision making. A popular approach to maximize the model's tail performance is to minimize the CVaR (Conditional Value at Risk) loss, which computes the average risk over the tails of the loss. However, for classification tasks where models are evaluated by the zero-one loss, we show that if the classifiers are deterministic, then the minimizer of the average zero-one loss also minimizes the CVaR zero-one loss, suggesting that CVaR loss minimization is not helpful without additional assumptions. We circumvent this negative result by minimizing the CVaR loss over randomized classifiers, for which the minimizers of the average zero-one loss and the CVaR zero-one loss are no longer the same, so minimizing the latter can lead to better tail performance. To learn such randomized classifiers, we propose the Boosted CVaR Classification framework which is motivated by a direct relationship between CVaR and a classical boosting algorithm called LPBoost. Based on this framework, we design an algorithm called \$\alpha\$-AdaLPBoost. We empirically evaluate our proposed algorithm on four benchmark datasets and show that it achieves higher tail performance than deterministic model training methods.

[Disentangled Contrastive Learning on Graphs](#)

- Haoyang Li, Xin Wang, Ziwei Zhang, Zehuan Yuan, Hang Li, Wenwu Zhu
- abstract@[open-review](#): Recently, self-supervised learning for graph neural networks (GNNs) has attracted considerable attention because of their notable successes in learning the representation of graph-structure data. However, the formation of a real-world graph typically arises from the highly complex interaction of many latent factors. The existing self-supervised learning methods for GNNs are inherently holistic and neglect the entanglement of the latent factors, resulting in the learned representations suboptimal for downstream tasks and difficult to be interpreted. Learning disentangled graph representations with self-supervised learning poses great challenges and remains largely ignored by the existing literature. In this paper, we introduce the Disentangled Graph Contrastive Learning (DGCL) method, which is able to learn disentangled graph-level representations with self-supervision. In particular, we first identify the latent factors of the input graph and derive its factorized representations. Each of the factorized representations describes a latent and disentangled aspect pertinent to a specific latent factor of the graph. Then we propose a novel factor-wise discrimination objective in a contrastive learning manner, which can force the factorized representations to independently reflect the expressive information from different latent factors. Extensive experiments on both synthetic and real-world datasets demonstrate the superiority of our method against several state-of-the-art baselines.

[Widening the Pipeline in Human-Guided Reinforcement Learning with Explanation and Context-Aware Data Augmentation](#)

- Lin Guan, Mudit Verma, Suna (Sihang) Guo, Ruohan Zhang, Subbarao Kambhampati
- abstract@[open-review](#): Human explanation (e.g., in terms of feature importance) has been recently used to extend the communication channel between human and agent in interactive machine learning. Under this setting, human trainers provide not only the ground truth but also some form of explanation. However, this kind of human guidance was only investigated in supervised learning tasks, and it remains unclear how to best incorporate this type of human knowledge into deep reinforcement learning. In this paper, we present the first study of using human visual explanations in human-in-the-loop reinforcement learning (HIRL). We focus on the task of learning from feedback, in which the human trainer not only gives binary evaluative "good" or "bad" feedback for queried state-action pairs, but also provides a visual explanation by annotating relevant features in images. We propose EXPAND (EXPlanation AugmeNted feedBack) to encourage the model to encode task-relevant features through a context-aware data augmentation that only perturbs irrelevant features in human salient information. We choose five tasks, namely Pixel-Taxi and four Atari games, to evaluate the performance and sample efficiency of this approach. We show that our method significantly outperforms methods leveraging human explanation that are adapted from supervised learning, and Human-in-the-loop RL baselines that only utilize evaluative feedback.

[SOLQ: Segmenting Objects by Learning Queries](#)

- Bin Dong, Fangao Zeng, Tiancai Wang, Xiangyu Zhang, Yichen Wei
- abstract@[open-review](#): In this paper, we propose an end-to-end framework for instance segmentation. Based on the recently introduced DETR, our method, termed SOLQ, segments objects by learning unified queries. In SOLQ, each query represents one object and has multiple representations: class, location and mask. The object queries learned perform classification, box regression and mask encoding simultaneously in an unified vector form. During training phase, the mask vectors encoded are supervised by the compression coding of raw spatial masks. In inference time, mask vectors produced can be directly transformed to spatial masks by the inverse process of compression coding. Experimental results show that SOLQ can achieve state-of-the-art performance, surpassing most of existing approaches. Moreover, the joint learning of unified query representation can greatly improve the detection performance of DETR. We hope our SOLQ can serve as a strong baseline for the Transformer-based instance segmentation.

[Extending Lagrangian and Hamiltonian Neural Networks with Differentiable Contact Models](#)

- Yaofeng Desmond Zhong, Biswadip Dey, Amit Chakraborty
- abstract@[open-review](#): The incorporation of appropriate inductive bias plays a critical role in learning dynamics from data. A growing body of work has been exploring ways to enforce energy conservation in the learned dynamics by encoding Lagrangian or Hamiltonian dynamics into the neural network architecture. These existing approaches are based on differential equations, which do not allow discontinuity in the states and thereby limit the class of systems one can learn. However, in reality, most physical systems, such as legged robots and robotic manipulators, involve contacts and collisions, which introduce discontinuities in the states. In this paper, we introduce a differentiable contact model, which can capture contact mechanics: frictionless/frictional, as well as elastic/inelastic. This model can also accommodate inequality constraints, such as limits on the joint angles. The proposed contact model extends the scope of Lagrangian and Hamiltonian neural networks by allowing simultaneous learning of contact and system properties. We demonstrate this framework on a series of challenging 2D and 3D physical systems with different coefficients of restitution and friction. The learned dynamics can be used as a differentiable physics simulator for downstream gradient-based optimization tasks, such as planning and control.

[Best-case lower bounds in online learning](#)

- Cristóbal Guzmán, Nishant Mehta, Ali Mortazavi
- abstract@[open-review](#): Much of the work in online learning focuses on the study of sublinear upper bounds on the regret. In this work, we initiate the study of best-case lower bounds in online convex optimization, wherein we bound the largest improvement an algorithm can obtain relative to the single best action in hindsight. This problem is motivated by the goal of better understanding the adaptivity of a learning algorithm. Another motivation comes from fairness: it is known that best-case lower bounds are instrumental in obtaining algorithms for decision-theoretic online learning (DTOL) that satisfy a notion of group fairness. Our contributions are a general method to provide best-case lower bounds in Follow The Regularized Leader (FTRL) algorithms with time-varying regularizers, which we use to show that best-case lower bounds are of the same order as existing upper regret bounds: this includes situations with a fixed learning rate, decreasing learning rates, timeless methods, and adaptive gradient methods. In stark contrast, we show that the linearized version of FTRL can attain negative linear regret. Finally, in DTOL with two experts and binary losses, we fully characterize the best-case sequences, which provides a finer understanding of the best-case lower bounds.

[A Comprehensively Tight Analysis of Gradient Descent for PCA](#)

- Zhiqiang Xu, Ping Li
- abstract@[open-review](#): We study the Riemannian gradient method for PCA on which a crucial fact is that despite the simplicity of the considered setting, i.e., deterministic version of Krasulina's method, the convergence rate has not been well-understood yet. In this work, we provide a general tight analysis for the gap-dependent rate at $\$O(\frac{1}{\Delta} \log \frac{1}{\epsilon})\$$ that holds for any real symmetric matrix. More importantly, when the gap Δ is significantly smaller than the target accuracy ϵ on the objective sub-optimality of the final solution, the rate of this type is actually not tight any more, which calls for a worst-case rate. We further give the first worst-case analysis that achieves a rate of convergence at $\$O(\frac{1}{\epsilon} \log \frac{1}{\epsilon})\$$. The two analyses naturally roll out a comprehensively tight convergence rate at $\$O(\frac{1}{\max\{\Delta, \epsilon\}})\$$. Particularly, our gap-dependent analysis suggests a new promising learning rate for stochastic variance reduced PCA algorithms. Experiments are conducted to confirm our findings as well.

[On Robust Optimal Transport: Computational Complexity and Barycenter Computation](#)

- Khang Le, Huy Nguyen, Quang M Nguyen, Tung Pham, Hung Bui, Nhat Ho
- abstract@[open-review](#): We consider robust variants of the standard optimal transport, named robust optimal transport, where marginal constraints are relaxed via Kullback-Leibler divergence. We show that Sinkhorn-based algorithms can approximate the optimal cost of robust optimal transport in $\widetilde{O}(\frac{n^2}{\epsilon})$ time, in which n is the number of supports of the probability distributions and ϵ is the desired error. Furthermore, we investigate a fixed-support robust barycenter problem between m discrete probability distributions with at most n number of supports and develop an approximating algorithm based on iterative Bregman projections (IBP). For the specific case $m = 2$, we show that this algorithm can approximate the optimal barycenter value in $\widetilde{O}(\frac{mn^2}{\epsilon})$ time, thus being better than the previous complexity $\widetilde{O}(\frac{mn^2}{\epsilon})$ of the IBP algorithm for approximating the Wasserstein barycenter.

[Asymptotically Best Causal Effect Identification with Multi-Armed Bandits](#)

- Alan Malek, Silvia Chiappa
- abstract@[open-review](#): This paper considers the problem of selecting a formula for identifying a causal quantity of interest among a set of available formulas. We assume an online setting in which the investigator may alter the data collection mechanism in a data-dependent way with the aim of identifying the formula with lowest asymptotic variance in as few samples as possible. We formalize this setting by using the best-arm-identification bandit framework where the standard goal of learning the arm with the lowest loss is replaced with the goal of learning the arm that will produce the best estimate. We introduce new tools for constructing finite-sample confidence bounds on estimates of the asymptotic variance that account for the estimation of potentially complex nuisance functions, and adapt the best-arm-identification algorithms of LUCB and Successive Elimination to use these bounds. We validate our method by providing upper bounds on the sample complexity and an empirical study on artificially generated data.

[Learning rule influences recurrent network representations but not attractor structure in decision-making tasks](#)

- Brandon McMahan, Michael Kleinman, Jonathan Kao
- abstract@[open-review](#): Recurrent neural networks (RNNs) are popular tools for studying computational dynamics in neurobiological circuits. However, due to the dizzying array of design choices, it is unclear if computational dynamics unearthed from RNNs provide reliable neurobiological inferences. Understanding the effects of design choices on RNN computation is valuable in two ways. First, invariant properties that persist in RNNs across a wide range of design choices are more likely to be candidate neurobiological mechanisms. Second, understanding what design choices lead to similar dynamical solutions reduces the burden of imposing that all design choices be totally faithful replications of biology. We focus our investigation on how RNN learning rule and task design affect RNN computation. We trained large populations of RNNs with different, but commonly used, learning rules on decision-making tasks inspired by neuroscience literature. For relatively complex tasks, we find that attractor topology is invariant to the choice of learning rule, but representational geometry is not. For simple tasks, we find that attractor topology depends on task input noise. However, when a task becomes increasingly complex, RNN attractor topology becomes invariant to input noise. Together, our results suggest that RNN dynamics are robust across learning rules but can be sensitive to the training task design, especially for simpler tasks.

[Few-Shot Segmentation via Cycle-Consistent Transformer](#)

- Gengwei Zhang, Guoliang Kang, Yi Yang, Yunchao Wei
- abstract@[open-review](#): Few-shot segmentation aims to train a segmentation model that can fast adapt to novel classes with few exemplars. The conventional training paradigm is to learn to make predictions on query images conditioned on the features from support images. Previous methods only utilized the semantic-level prototypes of support images as the conditional information. These methods cannot utilize all pixel-wise support information for the query predictions, which is however critical for the segmentation task. In this paper, we focus on utilizing pixel-wise relationships between support and target images to facilitate the few-shot semantic segmentation task. We design a novel Cycle-Consistent Transformer (CyCTR) module to aggregate pixel-wise support features into query ones. CyCTR performs cross-attention between features from different images, i.e. support and query images. We observe that there may exist unexpected irrelevant pixel-level support features. Directly performing cross-attention may aggregate these features from support to query and bias the query features. Thus, we propose using a novel cycle-consistent attention mechanism to filter out possible harmful support features and encourage query features to attend to the most informative pixels from support images. Experiments on all few-shot segmentation benchmarks demonstrate that our proposed CyCTR leads to remarkable improvement compared to previous state-of-the-art methods. Specifically, on Pascal-5^i and COCO-20^i datasets, we achieve 66.6% and 45.6% mIoU for 5-shot segmentation, outperforming previous state-of-the-art by 4.6% and 7.1% respectively.

[DropGNN: Random Dropouts Increase the Expressiveness of Graph Neural Networks](#)

- Pál András Papp, Karolis Martinkus, Lukas Faber, Roger Wattenhofer
- abstract@[open-review](#): This paper studies Dropout Graph Neural Networks (DropGNNs), a new approach that aims to overcome the limitations of standard GNN frameworks. In DropGNNs, we execute multiple runs of a GNN on the input graph, with some of the nodes randomly and independently dropped in each of these runs. Then, we combine the results of these runs to obtain the final result. We prove that DropGNNs can distinguish various graph neighborhoods that cannot be separated by message passing GNNs. We derive theoretical bounds for the number of runs required to ensure a reliable distribution of dropouts, and we prove several properties regarding the expressive capabilities and limits of DropGNNs. We experimentally validate our theoretical findings on expressiveness. Furthermore, we show that DropGNNs perform competitively on established GNN benchmarks.

[Photonic Differential Privacy with Direct Feedback Alignment](#)

- Ruben Ohana, Hamlet Medina, Julien Launay, Alessandro Cappelli, Iacopo Poli, Liva Ralaivola, Alain Rakotomamonjy
- abstract@[open-review](#): Optical Processing Units (OPUs) -- low-power photonic chips dedicated to large scale random projections -- have been used in previous work to train deep neural networks using Direct Feedback Alignment (DFA), an effective alternative to backpropagation. Here, we demonstrate how to leverage the intrinsic noise of optical random projections to build a differentially private DFA mechanism, making OPUs a solution of choice to provide a \emph{private-by-design} training. We provide a theoretical analysis of our adaptive privacy mechanism, carefully measuring how the noise of optical random projections propagates in the process and gives rise to provable Differential Privacy. Finally, we conduct experiments demonstrating the ability of our learning procedure to achieve solid end-task performance.

[Searching Parameterized AP Loss for Object Detection](#)

- Tao Chenxin, Zizhang Li, Xizhou Zhu, Gao Huang, Yong Liu, jifeng dai
- abstract@[open-review](#): Loss functions play an important role in training deep-network-based object detectors. The most widely used evaluation metric for object detection is Average Precision (AP), which captures the performance of localization and classification sub-tasks simultaneously. However, due to the non-differentiable nature of the AP metric, traditional object detectors adopt separate differentiable losses for the two sub-tasks. Such a mis-alignment issue may well lead to performance degradation. To address this, existing works seek to design surrogate losses for the AP metric manually, which requires expertise and may still be sub-optimal. In this paper, we propose Parameterized AP Loss, where parameterized functions are introduced to substitute the non-differentiable components in the AP calculation. Different AP approximations are thus represented by a family of parameterized functions in a unified formula. Automatic parameter search algorithm is then employed to search for the optimal parameters. Extensive experiments on the COCO benchmark with three different object detectors (i.e., RetinaNet, Faster R-CNN, and Deformable DETR) demonstrate that the proposed Parameterized AP Loss consistently outperforms existing handcrafted losses. Code shall be released.

[Fair Exploration via Axiomatic Bargaining](#)

- Jackie Baek, Vivek Farias
- abstract@[open-review](#): Motivated by the consideration of fairly sharing the cost of exploration between multiple groups in learning problems, we develop the Nash bargaining solution in the context of multi-armed bandits. Specifically, the 'grouped' bandit associated with any multi-armed bandit problem associates, with each time step, a single group from some finite set of groups. The utility gained by a given group under some learning policy is naturally viewed as the reduction in that group's regret relative to the regret that group would have incurred 'on its own'. We derive policies that yield the Nash bargaining solution relative to the set of incremental utilities possible under any policy. We show that on the one hand, the 'price of fairness' under such policies is limited, while on the other hand, regret optimal policies are arbitrarily unfair under generic conditions. Our theoretical development is complemented by a case study on contextual bandits for warfarin dosing where we are concerned with the cost of exploration across multiple races and age groups.

[Unifying lower bounds on prediction dimension of convex surrogates](#)

- Jessica Finocchiaro, Rafael Frongillo, Bo Waggoner
- abstract@[open-review](#): The convex consistency dimension of a supervised learning task is the lowest prediction dimension \$d\$ such that there exists a convex surrogate \$L : \mathbb{R}^d \times \mathcal{Y} \rightarrow \mathbb{R}\$ that is consistent for the given task. We present a new tool based on property elicitation, \$d\$-flats, for lower-bounding convex consistency dimension. This tool unifies approaches from a variety of domains, including continuous and

discrete prediction problems. We use $\$d\$$ -flats to obtain a new lower bound on the convex consistency dimension of risk measures, resolving an open question due to Frongillo and Kash (NeurIPS 2015). In discrete prediction settings, we show that the $\$d\$$ -flats approach recovers and even tightens previous lower bounds using feasible subspace dimension.

[Ultrahyperbolic Neural Networks](#)

- Marc Law
- abstract@[open-review](#): Riemannian space forms, such as the Euclidean space, sphere and hyperbolic space, are popular and powerful representation spaces in machine learning. For instance, hyperbolic geometry is appropriate to represent graphs without cycles and has been used to extend Graph Neural Networks. Recently, some pseudo-Riemannian space forms that generalize both hyperbolic and spherical geometries have been exploited to learn a specific type of nonparametric embedding called ultrahyperbolic. The lack of geodesic between every pair of ultrahyperbolic points makes the task of learning parametric models (e.g., neural networks) difficult. This paper introduces a method to learn parametric models in ultrahyperbolic space. We experimentally show the relevance of our approach in the tasks of graph and node classification.

[NeuroMLR: Robust & Reliable Route Recommendation on Road Networks](#)

- Jayant Jain, Vrittika Bagadia, Sahil Manchanda, Sayan Ranu
- abstract@[open-review](#): Predicting the most likely route from a source location to a destination is a core functionality in mapping services. Although the problem has been studied in the literature, two key limitations remain to be addressed. First, our study reveals that a significant portion of the routes recommended by existing methods fail to reach the destination. Second, existing techniques are transductive in nature; hence, they fail to recommend routes if unseen roads are encountered at inference time. In this paper, we address these limitations through an inductive algorithm called NeuroMLR. NeuroMLR learns a generative model from historical trajectories by conditioning on three explanatory factors: the current location, the destination, and real-time traffic conditions. The conditional distributions are learned through a novel combination of Lipschitz embedding with Graph Convolutional Networks (GCN) using historical trajectory data. Through in-depth experiments on real-world datasets, we establish that NeuroMLR imparts significant improvement in accuracy over the state of the art. More importantly, NeuroMLR generalizes dramatically better to unseen data and the recommended routes reach the destination with much higher likelihood than existing techniques.

[Risk Bounds and Calibration for a Smart Predict-then-Optimize Method](#)

- Heyuan Liu, Paul Grigas
- abstract@[open-review](#): The predict-then-optimize framework is fundamental in practical stochastic decision-making problems: first predict unknown parameters of an optimization model, then solve the problem using the predicted values. A natural loss function in this setting is defined by measuring the decision error induced by the predicted parameters, which was named the Smart Predict-then-Optimize (SPO) loss by Elmachtoub and Grigas [2021]. Since the SPO loss is typically nonconvex and possibly discontinuous, Elmachtoub and Grigas [2021] introduced a convex surrogate, called the SPO+ loss, that importantly accounts for the underlying structure of the optimization model. In this paper, we greatly expand upon the consistency results for the SPO+ loss provided by Elmachtoub and Grigas [2021]. We develop risk bounds and uniform calibration results for the SPO+ loss relative to the SPO loss, which provide a quantitative way to transfer the excess surrogate risk to excess true risk. By combining our risk bounds with generalization bounds, we show that the empirical minimizer of the SPO+ loss achieves low excess true risk with high probability. We first demonstrate these results in the case when the feasible region of the underlying optimization problem is a polyhedron, and then we show that the results can be strengthened substantially when the feasible region is a level set of a strongly convex function. We perform experiments to empirically demonstrate the strength of the SPO+ surrogate, as compared to standard $\|\cdot\|_1$ and squared $\|\cdot\|_2^2$ prediction error losses, on portfolio allocation and cost-sensitive multi-class classification problems.

[Three-dimensional spike localization and improved motion correction for Neuropixels recordings](#)

- Julien Boussard, Erdem Varol, Hyun Dong Lee, Nishchal Dethé, Liam Paninski
- abstract@[open-review](#): Neuropixels (NP) probes are dense linear multi-electrode arrays that have rapidly become essential tools for studying the electrophysiology of large neural populations. Unfortunately, a number of challenges remain in analyzing the large datasets output by these probes. Here we introduce several new methods for extracting useful spiking information from NP probes. First, we use a simple point neuron model, together with a neural-network denoiser, to efficiently map spikes detected on the probe into three-dimensional localizations. Previous methods localized spikes in two dimensions only; we show that the new localization approach is significantly more robust and provides an improved feature set for clustering spikes according to neural identity ("spike sorting"). Next, we apply a Poisson denoising method to the resulting three-dimensional point-cloud representation of the data, and show that the resulting 3D images can be accurately registered over time, leading to improved tracking of time-varying neural activity over the probe, and in turn, crisper estimates of neural clusters over time. The code to reproduce our results and an example neuropixels dataset is provided in the supplementary material.

[Semi-Supervised Semantic Segmentation via Adaptive Equalization Learning](#)

- Hanzhe Hu, Fangyun Wei, Han Hu, Qiwei Ye, Jinshi Cui, Liwei Wang
- abstract@[open-review](#): Due to the limited and even imbalanced data, semi-supervised semantic segmentation tends to have poor performance on some certain categories, e.g., tailed categories in Cityscapes dataset which exhibits a long-tailed label distribution. Existing approaches almost all neglect this problem, and treat categories equally. Some popular approaches such as consistency regularization or pseudo-labeling may even harm the learning of under-performing categories, that the predictions or pseudo labels of these categories could be too inaccurate to guide the learning on the unlabeled data. In this paper, we look into this problem, and propose a novel framework for semi-supervised semantic segmentation, named adaptive equalization learning (AEL). AEL adaptively balances the training of well and badly performed categories, with a confidence bank to dynamically track category-wise performance during training. The confidence bank is leveraged as an indicator to tilt training towards under-performing categories, instantiated in three strategies: 1) adaptive Copy-Paste and CutMix data augmentation approaches which give more chance for under-performing categories to be copied or cut; 2) an adaptive data sampling approach to encourage pixels from under-performing category to be sampled; 3) a simple yet effective re-weighting method to alleviate the training noise raised by pseudo-labeling. Experimentally, AEL outperforms the state-of-the-art methods by a large margin on the Cityscapes and Pascal VOC benchmarks under various data partition protocols. Code is available at <https://github.com/hzhupku/SemiSeg-AEL>.

[On the Bias-Variance-Cost Tradeoff of Stochastic Optimization](#)

- Yifan Hu, Xin Chen, Niao He
- abstract@[open-review](#): We consider stochastic optimization when one only has access to biased stochastic oracles of the objective, and obtaining stochastic gradients with low biases comes at high costs. This setting captures a variety of optimization paradigms widely used in machine learning, such as conditional stochastic optimization, bilevel optimization, and distributionally robust optimization. We examine a family of multi-level Monte Carlo (MLMC) gradient methods that exploit a delicate trade-off among the bias, the variance, and the oracle cost. We provide a systematic study of their convergences and total computation complexities for strongly convex, convex, and nonconvex objectives, and demonstrate their superiority over the naive biased stochastic gradient method. Moreover, when applied to conditional stochastic optimization, the MLMC gradient methods significantly improve the best-known sample complexity in the literature.

[Averaging on the Bures-Wasserstein manifold: dimension-free convergence of gradient descent](#)

- Jason Altschuler, Sinho Chewi, Patrik R Gerber, Austin Stromme
- abstract@[open-review](#): We study first-order optimization algorithms for computing the barycenter of Gaussian distributions with respect to the optimal transport metric. Although the objective is geodesically non-convex, Riemannian gradient descent empirically converges rapidly, in fact faster than off-the-shelf methods such as Euclidean gradient descent and SDP solvers. This stands in stark contrast to the best-known theoretical results, which depend exponentially on the dimension. In this work, we prove new geodesic convexity results which provide stronger control of the iterates, yielding a dimension-free convergence rate. Our techniques also enable the analysis of two related notions of averaging, the entropically-regularized barycenter and the geometric median, providing the first convergence guarantees for these problems.

[Reinforcement Learning in Newcomblike Environments](#)

- James Bell, Linda Linsefors, Caspar Oesterheld, Joar Skalse
- abstract@[open-review](#): Newcomblike decision problems have been studied extensively in the decision theory literature, but they have so far been largely absent in the reinforcement learning literature. In this paper we study value-based reinforcement learning algorithms in the Newcomblike setting, and answer some of the fundamental theoretical questions about the behaviour of such algorithms in these environments. We show that a value-based reinforcement learning agent cannot converge to a policy that is not \emph{ratifiable}, i.e., does not only choose actions that are optimal given that policy. This gives us a powerful tool for reasoning about the limit behaviour of agents -- for example, it lets us show that there are Newcomblike environments in which a reinforcement learning agent cannot converge to any optimal policy. We show that a ratifiable policy always exists in our setting, but that there are cases in which a reinforcement learning agent normally cannot converge to it (and hence cannot converge at all). We also prove several results about the possible limit behaviours of agents in cases where they do not converge to any policy.

[Comprehensive Knowledge Distillation with Causal Intervention](#)

- Xiang Deng, Zhongfei Zhang
- abstract@[open-review](#): Knowledge distillation (KD) addresses model compression by distilling knowledge from a large model (teacher) to a smaller one (student). The existing distillation approaches mainly focus on using different criteria to align the sample representations learned by the student and the teacher, while they fail to transfer the class representations. Good class representations can benefit the sample representation learning by shaping the sample representation distribution. On the other hand, the existing approaches enforce the student to fully imitate the teacher while ignoring the fact that the teacher is typically not perfect. Although the teacher has learned rich and powerful representations, it also contains unignorable bias knowledge which is usually induced by the context prior (e.g., background) in the training data. To address these two issues, in this paper, we propose comprehensive, interventional distillation (CID) that captures both sample and class representations from the teacher while removing the bias with causal intervention. Different from the existing literature that uses the softened logits of the teacher as the training targets, CID considers the softened logits as the context information of an image, which is further used to remove the biased knowledge based on causal inference. Keeping the good representations while removing the bad bias enables CID to have a better generalization ability on test data and a better transferability across different datasets against the existing state-of-the-art approaches, which is demonstrated by extensive experiments on several benchmark datasets.

[Reinforcement Learning with Latent Flow](#)

- Wenling Shang, Xiaofei Wang, Aravind Srinivas, Aravind Rajeswaran, Yang Gao, Pieter Abbeel, Misha Laskin
- abstract@[open-review](#): Temporal information is essential to learning effective policies with Reinforcement Learning (RL). However, current state-of-the-art RL algorithms either assume that such information is given as part of the state space or, when learning from pixels, use the simple heuristic of frame-stacking to implicitly capture temporal information present in the image observations. This heuristic is in contrast to the current paradigm in video classification architectures, which utilize explicit encodings of temporal information through methods such as optical flow and two-stream architectures to achieve state-of-the-art performance. Inspired by leading video classification architectures, we introduce the Flow of Latents for Reinforcement Learning (Flare), a network architecture for RL that explicitly encodes temporal information through latent vector differences. We show that Flare recovers optimal performance in state-based RL without explicit access to the state velocity, solely with positional state information. Flare is the most sample efficient model-free pixel-based RL algorithm on the DeepMind Control suite when evaluated on the 500k and 1M step benchmarks across 5 challenging control tasks, and, when used with Rainbow DQN, outperforms the competitive baseline on Atari games at 100M time step benchmark across 8 challenging games.

[Understanding How Encoder-Decoder Architectures Attend](#)

- Kyle Aitken, Vinay Ramasesh, Yuan Cao, Niru Maheswaranathan
- abstract@[open-review](#): Encoder-decoder networks with attention have proven to be a powerful way to solve many sequence-to-sequence tasks. In these networks, attention aligns encoder and decoder states and is often used for visualizing network behavior. However, the mechanisms used by networks to generate appropriate attention matrices are still mysterious. Moreover, how these mechanisms vary depending on the particular architecture used for the encoder and decoder (recurrent, feed-forward, etc.) are also not well understood. In this work, we investigate how encoder-decoder networks solve different sequence-to-sequence tasks. We introduce a way of decomposing hidden states over a sequence into temporal (independent of input) and input-driven (independent of sequence position) components. This reveals how attention matrices are formed: depending on the task requirements, networks rely more heavily on either the temporal or input-driven components. These findings hold across both recurrent and feed-forward architectures despite their differences in forming the temporal components. Overall, our results provide new insight into the inner workings of attention-based encoder-decoder networks.

[Latent Execution for Neural Program Synthesis Beyond Domain-Specific Languages](#)

- Xinyun Chen, Dawn Song, Yuandong Tian
- abstract@[open-review](#): Program synthesis from input-output (IO) examples has been a long-standing challenge. While recent works demonstrated limited success on domain-specific languages (DSL), it remains highly challenging to apply them to real-world programming languages, such as C. Due to complicated syntax and token variation, there are three major challenges: (1) unlike many DSLs, programs in languages like C need to compile first and are not executed via interpreters; (2) the program search space grows exponentially when the syntax and semantics of the programming language become more complex; and (3) collecting a large-scale dataset of real-world programs is non-trivial. As a first step to address these challenges, we propose LaSynth and show its efficacy in a restricted-C domain (i.e., C code with tens of tokens, with sequential, branching, loop and simple arithmetic operations but no library call). More specifically, LaSynth learns the latent representation to approximate the execution of partially generated programs, even if they are incomplete in syntax (addressing (1)). The learned execution significantly improves the performance of next token prediction over existing approaches, facilitating search (addressing (2)). Finally, once trained with randomly generated ground-truth programs and their IO pairs, LaSynth can synthesize more concise programs that resemble human-written code. Furthermore, retraining our model with these synthesized programs yields better performance with fewer samples for both Karel and C program synthesis, indicating the promise of leveraging the learned program synthesizer to improve the dataset quality for input-output program synthesis (addressing (3)). When evaluating on whether the program execution outputs match the IO pairs, LaSynth achieves 55.2% accuracy on generating simple C code with tens of tokens including loops and branches, outperforming existing approaches without executors by around 20%.

[Two steps to risk sensitivity](#)

- Christopher Gagne, Peter Dayan
- abstract@[open-review](#): Distributional reinforcement learning (RL) – in which agents learn about all the possible long-term consequences of their actions, and not just the expected value – is of great recent interest. One of the most important affordances of a distributional view is facilitating a modern, measured, approach to risk when outcomes are not completely certain. By contrast, psychological and neuroscientific investigations into decision making under risk have utilized a variety of more venerable theoretical models such as prospect theory that lack axiomatically desirable properties such as coherence. Here, we consider a particularly relevant risk measure for modeling human and animal planning, called conditional value-at-risk (CVaR), which quantifies worst-case outcomes (e.g., vehicle accidents or predation). We first adopt a conventional distributional approach to CVaR in a sequential setting and reanalyze the choices of human decision-makers in the well-known two-step task, revealing substantial risk aversion that had been lurking under stickiness and perseveration. We then consider a further critical property of risk sensitivity, namely time consistency, showing alternatives to this form of CVaR that enjoy this desirable characteristic. We use simulations to examine settings in which the various forms differ in ways that have implications for human and animal planning and behavior.

[DECAF: Generating Fair Synthetic Data Using Causally-Aware Generative Networks](#)

- Boris van Breugel, Trent Kyono, Jeroen Berrevoets, Mihaela van der Schaar
- abstract@[open-review](#): Machine learning models have been criticized for reflecting unfair biases in the training data. Instead of solving for this by introducing fair learning algorithms directly, we focus on generating fair synthetic data, such that any downstream learner is fair. Generating fair synthetic data from unfair data - while remaining truthful to the underlying data-generating process (DGP) - is non-trivial. In this paper, we introduce DECAF: a GAN-based fair synthetic data generator for tabular data. With DECAF we embed the DGP explicitly as a structural causal model in the input layers of the generator, allowing each variable to be reconstructed conditioned on its causal parents. This procedure enables inference time debiasing, where biased edges can be strategically removed for satisfying user-defined fairness requirements. The DECAF framework is versatile and compatible with several popular definitions of fairness. In our experiments, we show that DECAF successfully removes undesired bias and - in contrast to existing methods - is capable of generating high-quality synthetic data. Furthermore, we provide theoretical guarantees on the generator's convergence and the fairness of downstream models.

[EvoGrad: Efficient Gradient-Based Meta-Learning and Hyperparameter Optimization](#)

- Ondrej Bohdal, Yongxin Yang, Timothy Hospedales
- abstract@[open-review](#): Gradient-based meta-learning and hyperparameter optimization have seen significant progress recently, enabling practical end-to-end training of neural networks together with many hyperparameters. Nevertheless, existing approaches are relatively expensive as they need to compute second-order derivatives and store a longer computational graph. This cost prevents scaling them to larger network architectures. We present EvoGrad, a new approach to meta-learning that draws upon evolutionary techniques to more efficiently compute hypergradients. EvoGrad estimates hypergradient with respect to hyperparameters without calculating second-order gradients, or storing a longer computational graph, leading to significant improvements in efficiency. We evaluate EvoGrad on three substantial recent meta-learning applications, namely cross-domain few-shot learning with feature-wise transformations, noisy label learning with Meta-Weight-Net and low-resource cross-lingual learning with meta representation transformation. The results show that EvoGrad significantly improves efficiency and enables scaling meta-learning to bigger architectures such as from ResNet10 to ResNet34.

[Biological learning in key-value memory networks](#)

- Danil Tyulmankov, Ching Fang, Annapurna Vadaparty, Guangyu Robert Yang
- abstract@[open-review](#): In neuroscience, classical Hopfield networks are the standard biologically plausible model of long-term memory, relying on Hebbian plasticity for storage and attractor dynamics for recall. In contrast, memory-augmented neural networks in machine learning commonly use a key-value mechanism to store and read out memories in a single step. Such augmented networks achieve impressive feats of memory compared to traditional variants, yet their biological relevance is unclear. We propose an implementation of basic key-value memory that stores inputs using a combination of biologically plausible three-factor plasticity rules. The same rules are recovered when network parameters are meta-learned. Our network performs on par with classical Hopfield networks on autoassociative memory tasks and can be naturally extended to continual recall, heteroassociative memory, and sequence learning. Our results suggest a compelling alternative to the classical Hopfield network as a model of biological long-term memory.

[Correlated Stochastic Block Models: Exact Graph Matching with Applications to Recovering Communities](#)

- Miklos Racz, Anirudh Sridhar
- abstract@[open-review](#): We consider the task of learning latent community structure from multiple correlated networks. First, we study the problem of learning the latent vertex correspondence between two edge-correlated stochastic block models, focusing on the regime where the average degree is logarithmic in the number of vertices. We derive the precise information-theoretic threshold for exact recovery: above the threshold there exists an estimator that outputs the true correspondence with probability close to 1, while below it no estimator can recover the true correspondence with probability bounded away from 0. As an application of our results, we show how one can exactly recover the latent communities using \emph{multiple} correlated graphs in parameter regimes where it is information-theoretically impossible to do so using just a single graph.

[Twice regularized MDPs and the equivalence between robustness and regularization](#)

- Esther Derman, Matthieu Geist, Shie Mannor
- abstract@[open-review](#): Robust Markov decision processes (MDPs) aim to handle changing or partially known system dynamics. To solve them, one typically resorts to robust optimization methods. However, this significantly increases computational complexity and limits scalability in both learning and planning. On the other hand, regularized MDPs show more stability in policy learning without impairing time complexity. Yet, they generally do not encompass uncertainty in the model dynamics. In this work, we aim to learn robust MDPs using regularization. We first show that regularized MDPs are a particular instance of robust MDPs with uncertain reward. We thus establish that policy iteration on reward-robust MDPs can have the same time complexity as on regularized MDPs. We further extend this relationship to MDPs with uncertain transitions: this leads to a regularization term with an additional dependence on the value function. We finally generalize regularized MDPs to twice regularized MDPs ($R\gamma^2$ MDPs), i.e., MDPs with both value and policy regularization. The corresponding Bellman operators enable developing policy iteration schemes with convergence and robustness guarantees. It also reduces planning and learning in robust MDPs to regularized MDPs.

[Nearly Minimax Optimal Reinforcement Learning for Discounted MDPs](#)

- Jiafan He, Dongruo Zhou, Quanquan Gu
- abstract@[open-review](#): We study the reinforcement learning problem for discounted Markov Decision Processes (MDPs) under the tabular setting. We propose a model-based algorithm named UCBVI-\$\gamma\$, which is based on the \emph{optimism in the face of uncertainty principle} and the Bernstein-type bonus. We show that UCBVI-\$\gamma\$ achieves an $\tilde{O}(\sqrt{SAT}/((1-\gamma)^{1.5})\log(T))$ regret, where \$S\$ is the number of states, \$A\$ is the number of actions, \$\gamma\$ is the discount factor and \$T\$ is the number of steps. In addition, we construct a class of hard

MDPs and show that for any algorithm, the expected regret is at least $\tilde{\Omega}(\sqrt{SAT}/\{(1-\gamma)^{1.5}\})$. Our upper bound matches the minimax lower bound up to logarithmic factors, which suggests that UCBVI- γ is nearly minimax optimal for discounted MDPs.

[Sparse Deep Learning: A New Framework Immune to Local Traps and Miscalibration](#)

- Yan Sun, Wenjun Xiong, Faming Liang
- abstract@[open-review](#): Deep learning has powered recent successes of artificial intelligence (AI). However, the deep neural network, as the basic model of deep learning, has suffered from issues such as local traps and miscalibration. In this paper, we provide a new framework for sparse deep learning, which has the above issues addressed in a coherent way. In particular, we lay down a theoretical foundation for sparse deep learning and propose prior annealing algorithms for learning sparse neural networks. The former has successfully tamed the sparse deep neural network into the framework of statistical modeling, enabling prediction uncertainty correctly quantified. The latter can be asymptotically guaranteed to converge to the global optimum, enabling the validity of the down-stream statistical inference. Numerical result indicates the superiority of the proposed method compared to the existing ones.

[Calibrating Predictions to Decisions: A Novel Approach to Multi-Class Calibration](#)

- Shengjia Zhao, Michael Kim, Roshni Sahoo, Tengyu Ma, Stefano Ermon
- abstract@[open-review](#): When facing uncertainty, decision-makers want predictions they can trust. A machine learning provider can convey confidence to decision-makers by guaranteeing their predictions are distribution calibrated--- amongst the inputs that receive a predicted vector of class probabilities q , the actual distribution over classes is given by q . For multi-class prediction problems, however, directly optimizing predictions under distribution calibration tends to be infeasible, requiring sample complexity that grows exponentially in the number of classes C . In this work, we introduce a new notion---decision calibration---that requires the predicted distribution and true distribution over classes to be ``indistinguishable'' to downstream decision-makers. This perspective gives a new characterization of distribution calibration: a predictor is distribution calibrated if and only if it is decision calibrated with respect to all decision-makers. Our main result shows that under a mild restriction, unlike distribution calibration, decision calibration is actually feasible. We design a recalibration algorithm that provably achieves decision calibration efficiently, provided that the decision-makers have a bounded number of actions (e.g., polynomial in C). We validate our recalibration algorithm empirically: compared to existing methods, decision calibration improves decision-making on skin lesion and ImageNet classification with modern neural network predictors.

[Lower Bounds and Optimal Algorithms for Smooth and Strongly Convex Decentralized Optimization Over Time-Varying Networks](#)

- Dmitry Kovalev, Elnur Gasanov, Alexander Gasnikov, Peter Richtarik
- abstract@[open-review](#): We consider the task of minimizing the sum of smooth and strongly convex functions stored in a decentralized manner across the nodes of a communication network whose links are allowed to change in time. We solve two fundamental problems for this task. First, we establish $\{\epsilon_m$ the first lower bounds} on the number of decentralized communication rounds and the number of local computations required to find an ϵ -accurate solution. Second, we design two $\{\epsilon_m$ optimal algorithms} that attain these lower bounds: (i) a variant of the recently proposed algorithm ADOM (Kovalev et al, 2021) enhanced via a multi-consensus subroutine, which is optimal in the case when access to the dual gradients is assumed, and (ii) a novel algorithm, called ADOM+, which is optimal in the case when access to the primal gradients is assumed. We corroborate the theoretical efficiency of these algorithms by performing an experimental comparison with existing state-of-the-art methods.

[Testing Probabilistic Circuits](#)

- Yash Pralhad Pote, Kuldeep S Meel
- abstract@[open-review](#): Probabilistic circuits (PCs) are a powerful modeling framework for representing tractable probability distributions over combinatorial spaces. In machine learning and probabilistic programming, one is often interested in understanding whether the distributions learned using PCs are close to the desired distribution. Thus, given two probabilistic circuits, a fundamental problem of interest is to determine whether their distributions are close to each other. The primary contribution of this paper is a closeness test for PCs with respect to the total variation distance metric. Our algorithm utilizes two common PC queries, counting and sampling. In particular, we provide a poly-time probabilistic algorithm to check the closeness of two PCs, when the PCs support tractable approximate counting and sampling. We demonstrate the practical efficiency of our algorithmic framework via a detailed experimental evaluation of a prototype implementation against a set of 375 PC benchmarks. We find that our test correctly decides the closeness of all 375 PCs within 3600 seconds.

[Pseudo-Spherical Contrastive Divergence](#)

- Lantao Yu, Jiaming Song, Yang Song, Stefano Ermon
- abstract@[open-review](#): Energy-based models (EBMs) offer flexible distribution parametrization. However, due to the intractable partition function, they are typically trained via contrastive divergence for maximum likelihood estimation. In this paper, we propose pseudo-spherical contrastive divergence (PS-CD) to generalize maximum likelihood learning of EBMs. PS-CD is derived from the maximization of a family of strictly proper homogeneous scoring rules, which avoids the computation of the intractable partition function and provides a generalized family of learning objectives that include contrastive divergence as a special case. Moreover, PS-CD allows us to flexibly choose various learning objectives to train EBMs without additional computational cost or variational minimax optimization. Theoretical analysis on the proposed method and extensive experiments on both synthetic data and commonly used image datasets demonstrate the effectiveness and modeling flexibility of PS-CD, as well as its robustness to data contamination, thus showing its superiority over maximum likelihood and f-EBMs.

[NORESQA: A Framework for Speech Quality Assessment using Non-Matching References](#)

- Pranay Manocha, Buye Xu, Anurag Kumar
- abstract@[open-review](#): The perceptual task of speech quality assessment (SQA) is a challenging task for machines to do. Objective SQA methods that rely on the availability of the corresponding clean reference have been the primary go-to approaches for SQA. Clearly, these methods fail in real-world scenarios where the ground truth clean references are not available. In recent years, non-intrusive methods that train neural networks to predict ratings or scores have attracted much attention, but they suffer from several shortcomings such as lack of robustness, reliance on labeled data for training and so on. In this work, we propose a new direction for speech quality assessment. Inspired by human's innate ability to compare and assess the quality of speech signals even when they have non-matching contents, we propose a novel framework that predicts a subjective relative quality score for the given speech signal with respect to any provided reference without using any subjective data. We show that neural networks trained using our framework produce scores that correlate well with subjective mean opinion scores (MOS) and are also competitive to methods such as DNSMOS, which explicitly relies on MOS from humans for training networks. Moreover, our method also provides a natural way to embed quality-related information in neural networks, which we show is helpful for downstream tasks such as speech enhancement.

[AFEC: Active Forgetting of Negative Transfer in Continual Learning](#)

- Liyuan Wang, Mingtian Zhang, Zhongfan Jia, Qian Li, Chenglong Bao, Kaisheng Ma, Jun Zhu, Yi Zhong

- abstract@[open-review](#): Continual learning aims to learn a sequence of tasks from dynamic data distributions. Without accessing to the old training samples, knowledge transfer from the old tasks to each new task is difficult to determine, which might be either positive or negative. If the old knowledge interferes with the learning of a new task, i.e., the forward knowledge transfer is negative, then precisely remembering the old tasks will further aggravate the interference, thus decreasing the performance of continual learning. By contrast, biological neural networks can actively forget the old knowledge that conflicts with the learning of a new experience, through regulating the learning-triggered synaptic expansion and synaptic convergence. Inspired by the biological active forgetting, we propose to actively forget the old knowledge that limits the learning of new tasks to benefit continual learning. Under the framework of Bayesian continual learning, we develop a novel approach named Active Forgetting with synaptic Expansion-Convergence (AFEC). Our method dynamically expands parameters to learn each new task and then selectively combines them, which is formally consistent with the underlying mechanism of biological active forgetting. We extensively evaluate AFEC on a variety of continual learning benchmarks, including CIFAR-10 regression tasks, visual classification tasks and Atari reinforcement tasks, where AFEC effectively improves the learning of new tasks and achieves the state-of-the-art performance in a plug-and-play way.

[Heterogeneous Multi-player Multi-armed Bandits: Closing the Gap and Generalization](#)

- Chengshuai Shi, Wei Xiong, Cong Shen, Jing Yang
- abstract@[open-review](#): Despite the significant interests and many progresses in decentralized multi-player multi-armed bandits (MP-MAB) problems in recent years, the regret gap to the natural centralized lower bound in the heterogeneous MP-MAB setting remains open. In this paper, we propose BEACON -- Batched Exploration with Adaptive COmmunicatioN -- that closes this gap. BEACON accomplishes this goal with novel contributions in implicit communication and efficient exploration. For the former, we propose a novel adaptive differential communication (ADC) design that significantly improves the implicit communication efficiency. For the latter, a carefully crafted batched exploration scheme is developed to enable incorporation of the combinatorial upper confidence bound (CUCB) principle. We then generalize the existing linear-reward MP-MAB problems, where the system reward is always the sum of individually collected rewards, to a new MP-MAB problem where the system reward is a general (nonlinear) function of individual rewards. We extend BEACON to solve this problem and prove a logarithmic regret. BEACON bridges the algorithm design and regret analysis of combinatorial MAB (CMAB) and MP-MAB, two largely disjointed areas in MAB, and the results in this paper suggest that this previously ignored connection is worth further investigation.

[SWAD: Domain Generalization by Seeking Flat Minima](#)

- Junbum Cha, Sanghyuk Chun, Kyungjae Lee, Han-Cheol Cho, Seunghyun Park, Yunsung Lee, Sungrae Park
- abstract@[open-review](#): Domain generalization (DG) methods aim to achieve generalizability to an unseen target domain by using only training data from the source domains. Although a variety of DG methods have been proposed, a recent study shows that under a fair evaluation protocol, called DomainBed, the simple empirical risk minimization (ERM) approach works comparable to or even outperforms previous methods. Unfortunately, simply solving ERM on a complex, non-convex loss function can easily lead to sub-optimal generalizability by seeking sharp minima. In this paper, we theoretically show that finding flat minima results in a smaller domain generalization gap. We also propose a simple yet effective method, named Stochastic Weight Averaging Densely (SWAD), to find flat minima. SWAD finds flatter minima and suffers less from overfitting than does the vanilla SWA by a dense and overfit-aware stochastic weight sampling strategy. SWAD shows state-of-the-art performances on five DG benchmarks, namely PACS, VLCS, OfficeHome, TerraIncognita, and DomainNet, with consistent and large margins of +1.6% averagely on out-of-domain accuracy. We also compare SWAD with conventional generalization methods, such as data augmentation and consistency regularization methods, to verify that the remarkable performance improvements are originated from by seeking flat minima, not from better in-domain generalizability. Last but not least, SWAD is readily adaptable to existing DG methods without modification; the combination of SWAD and an existing DG method further improves DG performances. Source code is available at <https://github.com/khanrc/swad>.

[Autoformer: Decomposition Transformers with Auto-Correlation for Long-Term Series Forecasting](#)

- Haixu Wu, Jiehui Xu, Jianmin Wang, Mingsheng Long
- abstract@[open-review](#): Extending the forecasting time is a critical demand for real applications, such as extreme weather early warning and long-term energy consumption planning. This paper studies the long-term forecasting problem of time series. Prior Transformer-based models adopt various self-attention mechanisms to discover the long-range dependencies. However, intricate temporal patterns of the long-term future prohibit the model from finding reliable dependencies. Also, Transformers have to adopt the sparse versions of point-wise self-attentions for long series efficiency, resulting in the information utilization bottleneck. Going beyond Transformers, we design Autoformer as a novel decomposition architecture with an Auto-Correlation mechanism. We break with the pre-processing convention of series decomposition and renovate it as a basic inner block of deep models. This design empowers Autoformer with progressive decomposition capacities for complex time series. Further, inspired by the stochastic process theory, we design the Auto-Correlation mechanism based on the series periodicity, which conducts the dependencies discovery and representation aggregation at the sub-series level. Auto-Correlation outperforms self-attention in both efficiency and accuracy. In long-term forecasting, Autoformer yields state-of-the-art accuracy, with a 38% relative improvement on six benchmarks, covering five practical applications: energy, traffic, economics, weather and disease. Code is available at this repository: <https://github.com/thuml/Autoformer>.

[Predicting Event Memorability from Contextual Visual Semantics](#)

- Qianli Xu, Fen Fang, Ana Molino, Vigneshwaran Subbaraju, Joo-Hwee Lim
- abstract@[open-review](#): Episodic event memory is a key component of human cognition. Predicting event memorability, i.e., to what extent an event is recalled, is a tough challenge in memory research and has profound implications for artificial intelligence. In this study, we investigate factors that affect event memorability according to a cued recall process. Specifically, we explore whether event memorability is contingent on the event context, as well as the intrinsic visual attributes of image cues. We design a novel experiment protocol and conduct a large-scale experiment with 47 elder subjects over 3 months. Subjects' memory of life events is tested in a cued recall process. Using advanced visual analytics methods, we build a first-of-its-kind event memorability dataset (called R3) with rich information about event context and visual semantic features. Furthermore, we propose a contextual event memory network (CEMNet) that tackles multi-modal input to predict item-wise event memorability, which outperforms competitive benchmarks. The findings inform deeper understanding of episodic event memory, and open up a new avenue for prediction of human episodic memory. Source code is available at <https://github.com/ffzzy840304/Predicting-Event-Memorability>.

[Achieving Forgetting Prevention and Knowledge Transfer in Continual Learning](#)

- Zixuan Ke, Bing Liu, Nianzu Ma, Hu Xu, Lei Shu
- abstract@[open-review](#): Continual learning (CL) learns a sequence of tasks incrementally with the goal of achieving two main objectives: overcoming catastrophic forgetting (CF) and encouraging knowledge transfer (KT) across tasks. However, most existing techniques focus only on overcoming CF and have no mechanism to encourage KT, and thus do not do well in KT. Although several papers have tried to deal with both CF and KT, our experiments show that they suffer from serious CF when the tasks do not have much shared knowledge. Another observation is that most current CL methods do not use pre-trained models, but it has been shown that such models can significantly improve the end task performance. For example, in natural language processing, fine-tuning a BERT-like pre-trained language model is one of the most effective approaches. However, for CL, this approach suffers from serious CF. An interesting question is how to make the best use of pre-trained models for CL. This paper proposes a novel model called CTR to solve these problems. Our experimental results demonstrate the effectiveness of CTR

Bandits with many optimal arms

- Rianne de Heide, James Cheshire, Pierre Ménard, Alexandra Carpentier
- abstract@[open-review](#): We consider a stochastic bandit problem with a possibly infinite number of arms. We write p^* for the proportion of optimal arms and Δ for the minimal mean-gap between optimal and sub-optimal arms. We characterize the optimal learning rates both in the cumulative regret setting, and in the best-arm identification setting in terms of the problem parameters T (the budget), p^* and Δ . For the objective of minimizing the cumulative regret, we provide a lower bound of order $\Omega(\log(T)/(p^*\Delta))$ and a UCB-style algorithm with matching upper bound up to a factor of $\log(1/\Delta)$. Our algorithm needs p^* to calibrate its parameters, and we prove that this knowledge is necessary, since adapting to p^* in this setting is impossible. For best-arm identification we also provide a lower bound of order $\Omega(\exp(-cT\Delta^2p^*))$ on the probability of outputting a sub-optimal arm where $c > 0$ is an absolute constant. We also provide an elimination algorithm with an upper bound matching the lower bound up to a factor of order $\log(T)$ in the exponential, and that does not need p^* or Δ as parameter. Our results apply directly to the three related problems of competing against the j -th best arm, identifying an ϵ -good arm, and finding an arm with mean larger than a quantile of a known order.

Combiner: Full Attention Transformer with Sparse Computation Cost

- Hongyu Ren, Hanjun Dai, Zihang Dai, Mengjiao Yang, Jure Leskovec, Dale Schuurmans, Bo Dai
- abstract@[open-review](#): Transformers provide a class of expressive architectures that are extremely effective for sequence modeling. However, the key limitation of transformers is their quadratic memory and time complexity $\mathcal{O}(L^2)$ with respect to the sequence length in attention layers, which restricts application in extremely long sequences. Most existing approaches leverage sparsity or low-rank assumptions in the attention matrix to reduce cost, but sacrifice expressiveness. Instead, we propose Combiner, which provides full attention capability in each attention head while maintaining low computation and memory complexity. The key idea is to treat the self-attention mechanism as a conditional expectation over embeddings at each location, and approximate the conditional distribution with a structured factorization. Each location can attend to all other locations, either via direct attention, or through indirect attention to abstractions, which are again conditional expectations of embeddings from corresponding local regions. We show that most sparse attention patterns used in existing sparse transformers are able to inspire the design of such factorization for full attention, resulting in the same sub-quadratic cost ($\mathcal{O}(L \log(L))$ or $\mathcal{O}(L \sqrt{L})$). Combiner is a drop-in replacement for attention layers in existing transformers and can be easily implemented in common frameworks. An experimental evaluation on both autoregressive and bidirectional sequence tasks demonstrates the effectiveness of this approach, yielding state-of-the-art results on several image and text modeling tasks.

Geometry Processing with Neural Fields

- Guandao Yang, Serge Belongie, Bharath Hariharan, Vladlen Koltun
- abstract@[open-review](#): Most existing geometry processing algorithms use meshes as the default shape representation. Manipulating meshes, however, requires one to maintain high quality in the surface discretization. For example, changing the topology of a mesh usually requires additional procedures such as remeshing. This paper instead proposes the use of neural fields for geometry processing. Neural fields can compactly store complicated shapes without spatial discretization. Moreover, neural fields are infinitely differentiable, which allows them to be optimized for objectives that involve higher-order derivatives. This raises the question: can geometry processing be done entirely using neural fields? We introduce loss functions and architectures to show that some of the most challenging geometry processing tasks, such as deformation and filtering, can be done with neural fields. Experimental results show that our methods are on par with the well-established mesh-based methods without committing to a particular surface discretization. Code is available at <https://github.com/stevenygd/NFGP>.

Contextual Recommendations and Low-Regret Cutting-Plane Algorithms

- Sreenivas Gollapudi, Guru Guruganesh, Kostas Kollias, Pasin Manurangsi, Renato Leme, Jon Schneider
- abstract@[open-review](#): We consider the following variant of contextual linear bandits motivated by routing applications in navigational engines and recommendation systems. We wish to learn a hidden d -dimensional value w^* . Every round, we are presented with a subset $X_t \subseteq \mathbb{R}^d$ of possible actions. If we choose (i.e. recommend to the user) action x_t , we obtain utility $\langle x_t, w^* \rangle$ but only learn the identity of the best action $\arg\max_{x \in X_t} \langle x, w^* \rangle$. We design algorithms for this problem which achieve regret $\mathcal{O}(d \log T)$ and $\mathcal{O}(d \log d)$. To accomplish this, we design novel cutting-plane algorithms with low “regret” -- the total distance between the true point w^* and the hyperplanes the separation oracle returns. We also consider the variant where we are allowed to provide a list of several recommendations. In this variant, we give an algorithm with $\mathcal{O}(d^2 \log d)$ regret and list size $\mathcal{O}(d)$. Finally, we construct nearly tight algorithms for a weaker variant of this problem where the learner only learns the identity of an action that is better than the recommendation. Our results rely on new algorithmic techniques in convex geometry (including a variant of Steiner’s formula for the centroid of a convex set) which may be of independent interest.

Speech Separation Using an Asynchronous Fully Recurrent Convolutional Neural Network

- Xiaolin Hu, Kai Li, Weiyi Zhang, Yi Luo, Jean-Marie Lemercier, Timo Gerkmann
- abstract@[open-review](#): Recent advances in the design of neural network architectures, in particular those specialized in modeling sequences, have provided significant improvements in speech separation performance. In this work, we propose to use a bio-inspired architecture called Fully Recurrent Convolutional Neural Network (FRCNN) to solve the separation task. This model contains bottom-up, top-down and lateral connections to fuse information processed at various time-scales represented by stages. In contrast to the traditional approach updating stages in parallel, we propose to first update the stages one by one in the bottom-up direction, then fuse information from adjacent stages simultaneously and finally fuse information from all stages to the bottom stage together. Experiments showed that this asynchronous updating scheme achieved significantly better results with much fewer parameters than the traditional synchronous updating scheme on speech separation. In addition, the proposed model achieved competitive or better results with high efficiency as compared to other state-of-the-art approaches on two benchmark datasets.

Reinforcement Learning Enhanced Explainer for Graph Neural Networks

- Caihua Shan, Yifei Shen, Yao Zhang, Xiang Li, Dongsheng Li
- abstract@[open-review](#): Graph neural networks (GNNs) have recently emerged as revolutionary technologies for machine learning tasks on graphs. In GNNs, the graph structure is generally incorporated with node representation via the message passing scheme, making the explanation much more challenging. Given a trained GNN model, a GNN explainer aims to identify a most influential subgraph to interpret the prediction of an instance (e.g., a node or a graph), which is essentially a combinatorial optimization problem over graph. The existing works solve this problem by continuous relaxation or search-based heuristics. But they suffer from key issues such as violation of message passing and hand-crafted heuristics, leading to inferior interpretability. To address these issues, we propose a RL-enhanced GNN explainer, RG-Explainer, which consists of three main components: starting point selection, iterative graph generation and stopping criteria learning. RG-Explainer could construct a connected explanatory subgraph by sequentially adding nodes from the boundary of the current generated graph, which is consistent with the message passing scheme. Further, we design an effective seed locator to select the starting point, and learn stopping criteria to generate superior explanations. Extensive experiments on both synthetic and real datasets show that RG-Explainer outperforms state-of-the-art GNN explainers. Moreover, RG-Explainer can be applied in the inductive setting, demonstrating its better generalization ability.

NAS-Bench-x11 and the Power of Learning Curves

- Shen Yan, Colin White, Yash Savani, Frank Hutter
- abstract@[open-review](#): While early research in neural architecture search (NAS) required extreme computational resources, the recent releases of tabular and surrogate benchmarks have greatly increased the speed and reproducibility of NAS research. However, two of the most popular benchmarks do not provide the full training information for each architecture. As a result, on these benchmarks it is not possible to evaluate many types of multi-fidelity algorithms, such as learning curve extrapolation, that require evaluating architectures at arbitrary epochs. In this work, we present a method using singular value decomposition and noise modeling to create surrogate benchmarks, NAS-Bench-111, NAS-Bench-311, and NAS-Bench-NLP11, that output the full training information for each architecture, rather than just the final validation accuracy. We demonstrate the power of using the full training information by introducing a learning curve extrapolation framework to modify single-fidelity algorithms, showing that it leads to improvements over popular single-fidelity algorithms which claimed to be state-of-the-art upon release.

[Observation-Free Attacks on Stochastic Bandits](#)

- Yinglun Xu, Bhuvesh Kumar, Jacob D. Abernethy
- abstract@[open-review](#): We study data corruption attacks on stochastic multi arm bandit algorithms. Existing attack methodologies assume that the attacker can observe the multi arm bandit algorithm's realized behavior which is in contrast to the adversaries modeled in the robust multi arm bandit algorithms literature. To the best of our knowledge, we develop the first data corruption attack on stochastic multi arm bandit algorithms which works without observing the algorithm's realized behavior. Through this attack, we also discover a sufficient condition for a stochastic multi arm bandit algorithm to be susceptible to adversarial data corruptions. We show that any bandit algorithm that makes decisions just using the empirical mean reward, and the number of times that arm has been pulled in the past can suffer from linear regret under data corruption attacks. We further show that various popular stochastic multi arm bandit algorithms such UCB, ϵ -greedy and Thompson Sampling satisfy this sufficient condition and are thus prone to data corruption attacks. We further analyze the behavior of our attack for these algorithms and show that using only $O(T)$ corruptions, our attack can force these algorithms to select a potentially non-optimal target arm preferred by the attacker for all but $O(T)$ rounds.

[Learning Disentangled Behavior Embeddings](#)

- Changhao Shi, Sivan Schwartz, Shahar Levy, Shay Achvat, Maisan Abboud, Amir Ghanayim, Jackie Schiller, Gal Mishne
- abstract@[open-review](#): To understand the relationship between behavior and neural activity, experiments in neuroscience often include an animal performing a repeated behavior such as a motor task. Recent progress in computer vision and deep learning has shown great potential in the automated analysis of behavior by leveraging large and high-quality video datasets. In this paper, we design Disentangled Behavior Embedding (DBE) to learn robust behavioral embeddings from unlabeled, multi-view, high-resolution behavioral videos across different animals and multiple sessions. We further combine DBE with a stochastic temporal model to propose Variational Disentangled Behavior Embedding (VDBE), an end-to-end approach that learns meaningful discrete behavior representations and generates interpretable behavioral videos. Our models learn consistent behavior representations by explicitly disentangling the dynamic behavioral factors (pose) from time-invariant, non-behavioral nuisance factors (context) in a deep autoencoder, and exploit the temporal structures of pose dynamics. Compared to competing approaches, DBE and VDBE enjoy superior performance on downstream tasks such as fine-grained behavioral motif generation and behavior decoding.

[The Sensory Neuron as a Transformer: Permutation-Invariant Neural Networks for Reinforcement Learning](#)

- Yujin Tang, David Ha
- abstract@[open-review](#): In complex systems, we often observe complex global behavior emerge from a collection of agents interacting with each other in their environment, with each individual agent acting only on locally available information, without knowing the full picture. Such systems have inspired development of artificial intelligence algorithms in areas such as swarm optimization and cellular automata. Motivated by the emergence of collective behavior from complex cellular systems, we build systems that feed each sensory input from the environment into distinct, but identical neural networks, each with no fixed relationship with one another. We show that these sensory networks can be trained to integrate information received locally, and through communication via an attention mechanism, can collectively produce a globally coherent policy. Moreover, the system can still perform its task even if the ordering of its inputs is randomly permuted several times during an episode. These permutation invariant systems also display useful robustness and generalization properties that are broadly applicable. Interactive demo and videos of our results: <https://attentionneuron.github.io>

[Fast Extra Gradient Methods for Smooth Structured Nonconvex-Nonconcave Minimax Problems](#)

- Sucheol Lee, Donghwan Kim
- abstract@[open-review](#): Modern minimax problems, such as generative adversarial network and adversarial training, are often under a nonconvex-nonconcave setting, and developing an efficient method for such setting is of interest. Recently, two variants of the extragradient (EG) method are studied in that direction. First, a two-time-scale variant of the EG, named EG+, was proposed under a smooth structured nonconvex-nonconcave setting, with a slow $\mathcal{O}(1/k)$ rate on the squared gradient norm, where k denotes the number of iterations. Second, another variant of EG with an anchoring technique, named extra anchored gradient (EAG), was studied under a smooth convex-concave setting, yielding a fast $\mathcal{O}(1/k^2)$ rate on the squared gradient norm. Built upon EG+ and EAG, this paper proposes a two-time-scale EG with anchoring, named fast extragradient (FEG), that has a fast $\mathcal{O}(1/k^2)$ rate on the squared gradient norm for smooth structured nonconvex-nonconcave problems; the corresponding saddle-gradient operator satisfies the negative comonotonicity condition. This paper further develops its backtracking line-search version, named FEG-A, for the case where the problem parameters are not available. The stochastic analysis of FEG is also provided.

[Analysis of Sensing Spectral for Signal Recovery under a Generalized Linear Model](#)

- Junjie Ma, Ji Xu, Arian Maleki
- abstract@[open-review](#): We consider a nonlinear inverse problem $\mathbf{y} = f(\mathbf{A}\mathbf{x})$, where observations $\mathbf{y} \in \mathbb{R}^m$ are the componentwise nonlinear transformation of $\mathbf{A}\mathbf{x} \in \mathbb{R}^m$, $\mathbf{x} \in \mathbb{R}^n$ is the signal of interest and \mathbf{A} is a known linear mapping. By properly specifying the nonlinear processing function, this model can be particularized to many signal processing problems, including compressed sensing and phase retrieval. Our main goal in this paper is to understand the impact of sensing matrices, or more specifically the spectrum of sensing matrices, on the difficulty of recovering \mathbf{x} from \mathbf{y} . Towards this goal, we study the performance of one of the most successful recovery methods, i.e. the expectation propagation algorithm (EP). We define a notion for the spikiness of the spectrum of \mathbf{A} and show the importance of this measure in the performance of the EP. Whether the spikiness of the spectrum can hurt or help the recovery performance of EP depends on f . We define certain quantities based on the function f that enables us to describe the impact of the spikiness of the spectrum on EP recovery. Based on our framework, we are able to show that for instance, in phase-retrieval problems, matrices with spikier spectra are better for EP, while in 1-bit compressed sensing problems, less spiky (flatter) spectra offer better recoveries. Our results unify and substantially generalize the existing results that compare sub-Gaussian and orthogonal matrices, and provide a platform toward designing optimal sensing systems.

[Revisiting ResNets: Improved Training and Scaling Strategies](#)

- Irwan Bello, William Fedus, Xianzhi Du, Ekin Dogus Cubuk, Aravind Srinivas, Tsung-Yi Lin, Jonathon Shlens, Barret Zoph
- abstract@[open-review](#): Novel computer vision architectures monopolize the spotlight, but the impact of the model architecture is often conflated with simultaneous changes to training methodology and scaling strategies. Our work revisits the canonical ResNet and studies these three aspects in an effort to

disentangle them. Perhaps surprisingly, we find that training and scaling strategies may matter more than architectural changes, and further, that the resulting ResNets match recent state-of-the-art models. We show that the best performing scaling strategy depends on the training regime and offer two new scaling strategies: (1) scale model depth in regimes where overfitting can occur (width scaling is preferable otherwise); (2) increase image resolution more slowly than previously recommended. Using improved training and scaling strategies, we design a family of ResNet architectures, ResNet-RS, which are 1.7x - 2.7x faster than EfficientNets on TPUs, while achieving similar accuracies on ImageNet. In a large-scale semi-supervised learning setup, ResNet-RS achieves 86.2% top-1 ImageNet accuracy, while being 4.7x faster than EfficientNet-NoisyStudent. The training techniques improve transfer performance on a suite of downstream tasks (rivaling state-of-the-art self-supervised algorithms) and extend to video classification on Kinetics-400. We recommend practitioners use these simple revised ResNets as baselines for future research.

[Sparse Flows: Pruning Continuous-depth Models](#)

- Lucas Liebenwein, Ramin Hasani, Alexander Amini, Daniela Rus
- abstract@[open-review](#): Continuous deep learning architectures enable learning of flexible probabilistic models for predictive modeling as neural ordinary differential equations (ODEs), and for generative modeling as continuous normalizing flows. In this work, we design a framework to decipher the internal dynamics of these continuous depth models by pruning their network architectures. Our empirical results suggest that pruning improves generalization for neural ODEs in generative modeling. We empirically show that the improvement is because pruning helps avoid mode-collapse and flatten the loss surface. Moreover, pruning finds efficient neural ODE representations with up to 98% less parameters compared to the original network, without loss of accuracy. We hope our results will invigorate further research into the performance-size trade-offs of modern continuous-depth models.

[Spectrum-to-Kernel Translation for Accurate Blind Image Super-Resolution](#)

- Guangpin Tao, Xiaozhong Ji, Wenzhuo Wang, Shuo Chen, Chuming Lin, Yun Cao, Tong Lu, Donghao Luo, Ying Tai
- abstract@[open-review](#): Deep-learning based Super-Resolution (SR) methods have exhibited promising performance under non-blind setting where blur kernel is known; however, blur kernels of Low-Resolution (LR) images in different practical applications are usually unknown. It may lead to a significant performance drop when degradation process of training images deviates from that of real images. In this paper, we propose a novel blind SR framework to super-resolve LR images degraded by arbitrary blur kernel with accurate kernel estimation in frequency domain. To our best knowledge, this is the first deep learning method which conducts blur kernel estimation in frequency domain. Specifically, we first demonstrate that feature representation in frequency domain is more conducive for blur kernel reconstruction than in spatial domain. Next, we present a Spectrum-to-Kernel (S\$2\$K) network to estimate general blur kernels in diverse forms. We use a conditional GAN (CGAN) combined with SR-oriented optimization target to learn the end-to-end translation from degraded images' spectra to unknown kernels. Extensive experiments on both synthetic and real-world images demonstrate that our proposed method sufficiently reduces blur kernel estimation error, thus enables the off-the-shelf non-blind SR methods to work under blind setting effectively, and achieves superior performance over state-of-the-art blind SR methods, averagely by 1.39dB, 0.48dB (Gaussian kernels) and 6.15dB, 4.57dB (motion kernels) for scales \$2\times\$ and \$4\times\$ respectively.

[On the Rate of Convergence of Regularized Learning in Games: From Bandits and Uncertainty to Optimism and Beyond](#)

- Angeliki Giannou, Emmanouil-Vasileios Vlatakis-Gkaragkounis, Panayotis Mertikopoulos
- abstract@[open-review](#): In this paper, we examine the convergence rate of a wide range of regularized methods for learning in games. To that end, we propose a unified algorithmic template that we call “follow the generalized leader” (FTGL), and which includes as special cases the canonical “follow the regularized leader” algorithm, its optimistic variants, extra-gradient schemes, and many others. The proposed framework is also sufficiently flexible to account for several different feedback models – from full information to bandit feedback. In this general setting, we show that FTGL algorithms converge locally to strict Nash equilibria at a rate which does not depend on the level of uncertainty faced by the players, but only on the geometry of the regularizer near the equilibrium. In particular, we show that algorithms based on entropic regularization – like the exponential weights algorithm – enjoy a linear convergence rate, while Euclidean projection methods converge to equilibrium in a finite number of iterations, even with bandit feedback.

[SLAPS: Self-Supervision Improves Structure Learning for Graph Neural Networks](#)

- Bahare Fatemi, Layla El Asri, Seyed Mehran Kazemi
- abstract@[open-review](#): Graph neural networks (GNNs) work well when the graph structure is provided. However, this structure may not always be available in real-world applications. One solution to this problem is to infer a task-specific latent structure and then apply a GNN to the inferred graph. Unfortunately, the space of possible graph structures grows super-exponentially with the number of nodes and so the task-specific supervision may be insufficient for learning both the structure and the GNN parameters. In this work, we propose the Simultaneous Learning of Adjacency and GNN Parameters with Self-supervision, or SLAPS, a method that provides more supervision for inferring a graph structure through self-supervision. A comprehensive experimental study demonstrates that SLAPS scales to large graphs with hundreds of thousands of nodes and outperforms several models that have been proposed to learn a task-specific graph structure on established benchmarks.

[Aligning Pretraining for Detection via Object-Level Contrastive Learning](#)

- Fangyun Wei, Yue Gao, Zhirong Wu, Han Hu, Stephen Lin
- abstract@[open-review](#): Image-level contrastive representation learning has proven to be highly effective as a generic model for transfer learning. Such generality for transfer learning, however, sacrifices specificity if we are interested in a certain downstream task. We argue that this could be sub-optimal and thus advocate a design principle which encourages alignment between the self-supervised pretext task and the downstream task. In this paper, we follow this principle with a pretraining method specifically designed for the task of object detection. We attain alignment in the following three aspects: 1) object-level representations are introduced via selective search bounding boxes as object proposals; 2) the pretraining network architecture incorporates the same dedicated modules used in the detection pipeline (e.g. FPN); 3) the pretraining is equipped with object detection properties such as object-level translation invariance and scale invariance. Our method, called Selective Object COntastive learning (SoCo), achieves state-of-the-art results for transfer performance on COCO detection using a Mask R-CNN framework. Code is available at <https://github.com/hologerry/SoCo>.

[Double/Debiased Machine Learning for Dynamic Treatment Effects](#)

- Greg Lewis, Vasilis Syrgkanis
- abstract@[open-review](#): We consider the estimation of treatment effects in settings when multiple treatments are assigned over time and treatments can have a causal effect on future outcomes. We propose an extension of the double/debiased machine learning framework to estimate the dynamic effects of treatments and apply it to a concrete linear Markovian high-dimensional state space model and to general structural nested mean models. Our method allows the use of arbitrary machine learning methods to control for the high dimensional state, subject to a mean square error guarantee, while still allowing parametric estimation and construction of confidence intervals for the dynamic treatment effect parameters of interest. Our method is based on a sequential regression peeling process, which we show can be equivalently interpreted as a Neyman orthogonal moment estimator. This allows us to show root-n asymptotic normality of the estimated causal effects.

[Local Disentanglement in Variational Auto-Encoders Using Jacobian \$\|L\|_1\$ Regularization](#)

- Travers Rhodes, Daniel Lee
- abstract@[open-review](#): There have been many recent advances in representation learning; however, unsupervised representation learning can still struggle with model identification issues related to rotations of the latent space. Variational Auto-Encoders (VAEs) and their extensions such as β -VAEs have been shown to improve local alignment of latent variables with PCA directions, which can help to improve model disentanglement under some conditions. Borrowing inspiration from Independent Component Analysis (ICA) and sparse coding, we propose applying an L_1 loss to the VAE's generative Jacobian during training to encourage local latent variable alignment with independent factors of variation in images of multiple objects or images with multiple parts. We demonstrate our results on a variety of datasets, giving qualitative and quantitative results using information theoretic and modularity measures that show our added L_1 cost encourages local axis alignment of the latent representation with individual factors of variation.

[Design of Experiments for Stochastic Contextual Linear Bandits](#)

- Andrea Zanette, Kefan Dong, Jonathan N Lee, Emma Brunskill
- abstract@[open-review](#): In the stochastic linear contextual bandit setting there exist several minimax procedures for exploration with policies that are reactive to the data being acquired. In practice, there can be a significant engineering overhead to deploy these algorithms, especially when the dataset is collected in a distributed fashion or when a human in the loop is needed to implement a different policy. Exploring with a single non-reactive policy is beneficial in such cases. Assuming some batch contexts are available, we design a single stochastic policy to collect a good dataset from which a near-optimal policy can be extracted. We present a theoretical analysis as well as numerical experiments on both synthetic and real-world datasets.

[Encoding Spatial Distribution of Convolutional Features for Texture Representation](#)

- Yong Xu, Feng Li, Zhile Chen, Jinxiu Liang, Yuhui Quan
- abstract@[open-review](#): Existing convolutional neural networks (CNNs) often use global average pooling (GAP) to aggregate feature maps into a single representation. However, GAP cannot well characterize complex distributive patterns of spatial features while such patterns play an important role in texture-oriented applications, e.g., material recognition and ground terrain classification. In the context of texture representation, this paper addressed the issue by proposing Fractal Encoding (FE), a feature encoding module grounded by multi-fractal geometry. Considering a CNN feature map as a union of level sets of points lying in the 2D space, FE characterizes their spatial layout via a local-global hierarchical fractal analysis which examines the multi-scale power behavior on each level set. This enables a CNN to encode the regularity on the spatial arrangement of image features, leading to a robust yet discriminative spectrum descriptor. In addition, FE has trainable parameters for data adaptivity and can be easily incorporated into existing CNNs for end-to-end training. We applied FE to ResNet-based texture classification and retrieval, and demonstrated its effectiveness on several benchmark datasets.

[Training Certifiably Robust Neural Networks with Efficient Local Lipschitz Bounds](#)

- Yujia Huang, Huan Zhang, Yuanyuan Shi, J. Zico Kolter, Anima Anandkumar
- abstract@[open-review](#): Certified robustness is a desirable property for deep neural networks in safety-critical applications, and popular training algorithms can certify robustness of a neural network by computing a global bound on its Lipschitz constant. However, such a bound is often loose: it tends to over-regularize the neural network and degrade its natural accuracy. A tighter Lipschitz bound may provide a better tradeoff between natural and certified accuracy, but is generally hard to compute exactly due to non-convexity of the network. In this work, we propose an efficient and trainable \emph{local} Lipschitz upper bound by considering the interactions between activation functions (e.g. ReLU) and weight matrices. Specifically, when computing the induced norm of a weight matrix, we eliminate the corresponding rows and columns where the activation function is guaranteed to be a constant in the neighborhood of each given data point, which provides a provably tighter bound than the global Lipschitz constant of the neural network. Our method can be used as a plug-in module to tighten the Lipschitz bound in many certifiable training algorithms. Furthermore, we propose to clip activation functions (e.g., ReLU and MaxMin) with a learnable upper threshold and a sparsity loss to assist the network to achieve an even tighter local Lipschitz bound. Experimentally, we show that our method consistently outperforms state-of-the-art methods in both clean and certified accuracy on MNIST, CIFAR-10 and TinyImageNet datasets with various network architectures.

[Average-Reward Learning and Planning with Options](#)

- Yi Wan, Abhishek Naik, Rich Sutton
- abstract@[open-review](#): We extend the options framework for temporal abstraction in reinforcement learning from discounted Markov decision processes (MDPs) to average-reward MDPs. Our contributions include general convergent off-policy inter-option learning algorithms, intra-option algorithms for learning values and models, as well as sample-based planning variants of our learning algorithms. Our algorithms and convergence proofs extend those recently developed by Wan, Naik, and Sutton. We also extend the notion of option-interrupting behaviour from the discounted to the average-reward formulation. We show the efficacy of the proposed algorithms with experiments on a continuing version of the Four-Room domain.

[SSAL: Synergizing between Self-Training and Adversarial Learning for Domain Adaptive Object Detection](#)

- Muhammad Akhtar Munir, Muhammad Haris Khan, M. Sarfraz, Mohsen Ali
- abstract@[open-review](#): We study adapting trained object detectors to unseen domains manifesting significant variations of object appearance, viewpoints and backgrounds. Most current methods align domains by either using image or instance-level feature alignment in an adversarial fashion. This often suffers due to the presence of unwanted background and as such lacks class-specific alignment. A common remedy to promote class-level alignment is to use high confidence predictions on the unlabelled domain as pseudo labels. These high confidence predictions are often fallacious since the model is poorly calibrated under domain shift. In this paper, we propose to leverage model's predictive uncertainty to strike the right balance between adversarial feature alignment and class-level alignment. Specifically, we measure predictive uncertainty on class assignments and the bounding box predictions. Model predictions with low uncertainty are used to generate pseudo-labels for self-supervision, whereas the ones with higher uncertainty are used to generate tiles for an adversarial feature alignment stage. This synergy between tiling around the uncertain object regions and generating pseudo-labels from highly certain object regions allows us to capture both the image and instance level context during the model adaptation stage. We perform extensive experiments covering various domain shift scenarios. Our approach improves upon existing state-of-the-art methods with visible margins.

[Counterexample Guided RL Policy Refinement Using Bayesian Optimization](#)

- Briti Gangopadhyay, Pallab Dasgupta
- abstract@[open-review](#): Constructing Reinforcement Learning (RL) policies that adhere to safety requirements is an emerging field of study. RL agents learn via trial and error with an objective to optimize a reward signal. Often policies that are designed to accumulate rewards do not satisfy safety specifications. We present a methodology for counterexample guided refinement of a trained RL policy against a given safety specification. Our approach has two main components. The first component is an approach to discover failure trajectories using Bayesian optimization over multiple parameters of uncertainty from a policy learnt in a model-free setting. The second component selectively modifies the failure points of the policy using gradient-based updates. The approach has been tested on several RL environments, and we demonstrate that the policy can be made to respect the safety specifications through such targeted changes.

[Stable, Fast and Accurate: Kernelized Attention with Relative Positional Encoding](#)

- Shengjie Luo, Shanda Li, Tianle Cai, Di He, Dinglan Peng, Shuxin Zheng, Guolin Ke, Liwei Wang, Tie-Yan Liu
- abstract@[open-review](#): The attention module, which is a crucial component in Transformer, cannot scale efficiently to long sequences due to its quadratic complexity. Many works focus on approximating the dot-then-exponentiate softmax function in the original attention, leading to sub-quadratic or even linear-complexity Transformer architectures. However, we show that these methods cannot be applied to more powerful attention modules that go beyond the dot-then-exponentiate style, e.g., Transformers with relative positional encoding (RPE). Since in many state-of-the-art models, relative positional encoding is used as default, designing efficient Transformers that can incorporate RPE is appealing. In this paper, we propose a novel way to accelerate attention calculation for Transformers with RPE on top of the kernelized attention. Based upon the observation that relative positional encoding forms a Toeplitz matrix, we mathematically show that kernelized attention with RPE can be calculated efficiently using Fast Fourier Transform (FFT). With FFT, our method achieves $\mathcal{O}(n \log n)$ time complexity. Interestingly, we further demonstrate that properly using relative positional encoding can mitigate the training instability problem of vanilla kernelized attention. On a wide range of tasks, we empirically show that our models can be trained from scratch without any optimization issues. The learned model performs better than many efficient Transformer variants and is faster than standard Transformer in the long-sequence regime.

[Learning in Non-Cooperative Configurable Markov Decision Processes](#)

- Giorgia Ramponi, Alberto Maria Metelli, Alessandro Concetti, Marcello Restelli
- abstract@[open-review](#): The Configurable Markov Decision Process framework includes two entities: a Reinforcement Learning agent and a configurator that can modify some environmental parameters to improve the agent's performance. This presupposes that the two actors have the same reward functions. What if the configurator does not have the same intentions as the agent? This paper introduces the Non-Cooperative Configurable Markov Decision Process, a setting that allows having two (possibly different) reward functions for the configurator and the agent. Then, we consider an online learning problem, where the configurator has to find the best among a finite set of possible configurations. We propose two learning algorithms to minimize the configurator's expected regret, which exploits the problem's structure, depending on the agent's feedback. While a naive application of the UCB algorithm yields a regret that grows indefinitely over time, we show that our approach suffers only bounded regret. Furthermore, we empirically show the performance of our algorithm in simulated domains.

[Identification of Partially Observed Linear Causal Models: Graphical Conditions for the Non-Gaussian and Heterogeneous Cases](#)

- Jeffrey Adams, Niels Hansen, Kun Zhang
- abstract@[open-review](#): In causal discovery, linear non-Gaussian acyclic models (LiNGAMs) have been studied extensively. While the causally sufficient case is well understood, in many real problems the observed variables are not causally related. Rather, they are generated by latent variables, such as confounders and mediators, which may themselves be causally related. Existing results on the identification of the causal structure among the latent variables often require very strong graphical assumptions. In this paper, we consider partially observed linear models with either non-Gaussian or heterogeneous errors. In that case we give two graphical conditions which are necessary for identification of the causal structure. These conditions are closely related to sparsity of the causal edges. Together with one additional condition on the coefficients, which holds generically for any graph, the two graphical conditions are also sufficient for identifiability. These new conditions can be satisfied even when there is a large number of latent variables. We demonstrate the validity of our results on synthetic data.

[DIB-R++: Learning to Predict Lighting and Material with a Hybrid Differentiable Renderer](#)

- Wenzheng Chen, Joey Litalien, Jun Gao, Zian Wang, Clement Fuji Tsang, Sameh Khamis, Or Litany, Sanja Fidler
- abstract@[open-review](#): We consider the challenging problem of predicting intrinsic object properties from a single image by exploiting differentiable renderers. Many previous learning-based approaches for inverse graphics adopt rasterization-based renderers and assume naive lighting and material models, which often fail to account for non-Lambertian, specular reflections commonly observed in the wild. In this work, we propose DIBR++, a hybrid differentiable renderer which supports these photorealistic effects by combining rasterization and ray-tracing, taking the advantage of their respective strengths---speed and realism. Our renderer incorporates environmental lighting and spatially-varying material models to efficiently approximate light transport, either through direct estimation or via spherical basis functions. Compared to more advanced physics-based differentiable renderers leveraging path tracing, DIBR++ is highly performant due to its compact and expressive shading model, which enables easy integration with learning frameworks for geometry, reflectance and lighting prediction from a single image without requiring any ground-truth. We experimentally demonstrate that our approach achieves superior material and lighting disentanglement on synthetic and real data compared to existing rasterization-based approaches and showcase several artistic applications including material editing and relighting.

[Coresets for Time Series Clustering](#)

- Lingxiao Huang, K Sudhir, Nisheeth Vishnoi
- abstract@[open-review](#): We study the problem of constructing coresets for clustering problems with time series data. This problem has gained importance across many fields including biology, medicine, and economics due to the proliferation of sensors facilitating real-time measurement and rapid drop in storage costs. In particular, we consider the setting where the time series data on N entities is generated from a Gaussian mixture model with autocorrelations over k clusters in R^d . Our main contribution is an algorithm to construct coresets for the maximum likelihood objective for this mixture model. Our algorithm is efficient, and under a mild boundedness assumption on the covariance matrices of the underlying Gaussians, the size of the coreset is independent of the number of entities N and the number of observations for each entity, and depends only polynomially on k , d and ε , where ε is the error parameter. We empirically assess the performance of our coreset with synthetic data.

[A Variational Perspective on Diffusion-Based Generative Models and Score Matching](#)

- Chin-Wei Huang, Jae Hyun Lim, Aaron C. Courville
- abstract@[open-review](#): Discrete-time diffusion-based generative models and score matching methods have shown promising results in modeling high-dimensional image data. Recently, Song et al. (2021) show that diffusion processes that transform data into noise can be reversed via learning the score function, i.e. the gradient of the log-density of the perturbed data. They propose to plug the learned score function into an inverse formula to define a generative diffusion process. Despite the empirical success, a theoretical underpinning of this procedure is still lacking. In this work, we approach the (continuous-time) generative diffusion directly and derive a variational framework for likelihood estimation, which includes continuous-time normalizing flows as a special case, and can be seen as an infinitely deep variational autoencoder. Under this framework, we show that minimizing the score-matching loss is equivalent to maximizing a lower bound of the likelihood of the plug-in reverse SDE proposed by Song et al. (2021), bridging the theoretical gap.

[Online Active Learning with Surrogate Loss Functions](#)

- Giulia DeSalvo, Claudio Gentile, Tobias Sommer Thune
- abstract@[open-review](#): We derive a novel active learning algorithm in the streaming setting for binary classification tasks. The algorithm leverages weak labels to minimize the number of label requests, and trains a model to optimize a surrogate loss on a resulting set of labeled and weak-labeled points. Our algorithm jointly admits two crucial properties: theoretical guarantees in the general agnostic setting and a strong empirical performance. Our theoretical analysis shows that the algorithm attains favorable generalization and label complexity bounds, while our empirical study on 18 real-world datasets

demonstrate that the algorithm outperforms standard baselines, including the Margin Algorithm, or Uncertainty Sampling, a high-performing active learning algorithm favored by practitioners.

[Does Preprocessing Help Training Over-parameterized Neural Networks?](#)

- Zhao Song, Shuo Yang, Ruizhe Zhang
- abstract@[open-review](#): Deep neural networks have achieved impressive performance in many areas. Designing a fast and provable method for training neural networks is a fundamental question in machine learning. The classical training method requires paying $\Omega(mnd)$ cost for both forward computation and backward computation, where m is the width of the neural network, and we are given n training points in d -dimensional space. In this paper, we propose two novel preprocessing ideas to bypass this $\Omega(mnd)$ barrier: *First, by preprocessing the initial weights of the neural networks, we can train the neural network in $\widetilde{O}(m^{1-\Theta(1/d)} n d)$ cost per iteration.* Second, by preprocessing the input data points, we can train neural network in $\widetilde{O}(m^{4/5} nd)$ cost per iteration. From the technical perspective, our result is a sophisticated combination of tools in different fields, greedy-type convergence analysis in optimization, sparsity observation in practical work, high-dimensional geometric search in data structure, concentration and anti-concentration in probability. Our results also provide theoretical insights for a large number of previously established fast training methods. In addition, our classical algorithm can be generalized to the Quantum computation model. Interestingly, we can get a similar sublinear cost per iteration but avoid preprocessing initial weights or input data points.

[Causal Influence Detection for Improving Efficiency in Reinforcement Learning](#)

- Maximilian Seitzer, Bernhard Schölkopf, Georg Martius
- abstract@[open-review](#): Many reinforcement learning (RL) environments consist of independent entities that interact sparsely. In such environments, RL agents have only limited influence over other entities in any particular situation. Our idea in this work is that learning can be efficiently guided by knowing when and what the agent can influence with its actions. To achieve this, we introduce a measure of situation-dependent causal influence based on conditional mutual information and show that it can reliably detect states of influence. We then propose several ways to integrate this measure into RL algorithms to improve exploration and off-policy learning. All modified algorithms show strong increases in data efficiency on robotic manipulation tasks.

[LADA: Look-Ahead Data Acquisition via Augmentation for Deep Active Learning](#)

- Yoon-Yeong Kim, Kyungwoo Song, JoonHo Jang, Il-chul Moon
- abstract@[open-review](#): Active learning effectively collects data instances for training deep learning models when the labeled dataset is limited and the annotation cost is high. Data augmentation is another effective technique to enlarge the limited amount of labeled instances. The scarcity of labeled dataset leads us to consider the integration of data augmentation and active learning. One possible approach is a pipelined combination, which selects informative instances via the acquisition function and generates virtual instances from the selected instances via augmentation. However, this pipelined approach would not guarantee the informativeness of the virtual instances. This paper proposes Look-Ahead Data Acquisition via augmentation, or LADA framework, that looks ahead the effect of data augmentation in the process of acquisition. LADA jointly considers both 1) unlabeled data instance to be selected and 2) virtual data instance to be generated by data augmentation, to construct the acquisition function. Moreover, to generate maximally informative virtual instances, LADA optimizes the data augmentation policy to maximize the predictive acquisition score, resulting in the proposal of InfoSTN and InfoMixup. The experimental results of LADA show a significant improvement over the recent augmentation and acquisition baselines that were independently applied.

[Policy Optimization in Adversarial MDPs: Improved Exploration via Dilated Bonuses](#)

- Haipeng Luo, Chen-Yu Wei, Chung-Wei Lee
- abstract@[open-review](#): Policy optimization is a widely-used method in reinforcement learning. Due to its local-search nature, however, theoretical guarantees on global optimality often rely on extra assumptions on the Markov Decision Processes (MDPs) that bypass the challenge of global exploration. To eliminate the need of such assumptions, in this work, we develop a general solution that adds dilated bonuses to the policy update to facilitate global exploration. To showcase the power and generality of this technique, we apply it to several episodic MDP settings with adversarial losses and bandit feedback, improving and generalizing the state-of-the-art. Specifically, in the tabular case, we obtain $\widetilde{\mathcal{O}}(\sqrt{T})$ regret where T is the number of episodes, improving the $\widetilde{\mathcal{O}}(T^{2/3})$ regret bound by Shani et al. [2020]. When the number of states is infinite, under the assumption that the state-action values are linear in some low-dimensional features, we obtain $\widetilde{\mathcal{O}}(T^{2/3})$ regret with the help of a simulator, matching the result of Neu and Olkhovskaya [2020] while importantly removing the need of an exploratory policy that their algorithm requires. To our knowledge, this is the first algorithm with sublinear regret for linear function approximation with adversarial losses, bandit feedback, and no exploratory assumptions. Finally, we also discuss how to further improve the regret or remove the need of a simulator using dilated bonuses, when an exploratory policy is available.

[Multiclass versus Binary Differentially Private PAC Learning](#)

- Satchit Sivakumar, Mark Bun, Marco Gaboardi
- abstract@[open-review](#): We show a generic reduction from multiclass differentially private PAC learning to binary private PAC learning. We apply this transformation to a recently proposed binary private PAC learner to obtain a private multiclass learner with sample complexity that has a polynomial dependence on the multiclass Littlestone dimension and a poly-logarithmic dependence on the number of classes. This yields a doubly exponential improvement in the dependence on both parameters over learners from previous work. Our proof extends the notion of Ψ -dimension defined in work of Ben-David et al. [JCSS, 1995] to the online setting and explores its general properties.

[Adversarially Robust Change Point Detection](#)

- Mengchu Li, Yi Yu
- abstract@[open-review](#): Change point detection is becoming increasingly popular in many application areas. On one hand, most of the theoretically-justified methods are investigated in an ideal setting without model violations, or merely robust against identical heavy-tailed noise distribution across time and/or against isolate outliers; on the other hand, we are aware that there have been exponentially growing attacks from adversaries, who may pose systematic contamination on data to purposely create spurious change points or disguise true change points. In light of the timely need for a change point detection method that is robust against adversaries, we start with, arguably, the simplest univariate mean change point detection problem. The adversarial attacks are formulated through the Huber ε -contamination framework, which in particular allows the contamination distributions to be different at each time point. In this paper, we demonstrate a phase transition phenomenon in change point detection. This detection boundary is a function of the contamination proportion ε and is the first time shown in the literature. In addition, we derive the minimax-rate optimal localisation error rate, quantifying the cost of accuracy in terms of the contamination proportion. We propose a computationally feasible method, matching the minimax lower bound under certain conditions, saving for logarithmic factors. Extensive numerical experiments are conducted with comparisons to robust change point detection methods in the existing literature.

[Cycle Self-Training for Domain Adaptation](#)

- Hong Liu, Jianmin Wang, Mingsheng Long
- abstract@[open-review](#): Mainstream approaches for unsupervised domain adaptation (UDA) learn domain-invariant representations to narrow the domain shift, which are empirically effective but theoretically challenged by the hardness or impossibility theorems. Recently, self-training has been gaining momentum in UDA, which exploits unlabeled target data by training with target pseudo-labels. However, as corroborated in this work, under distributional shift, the pseudo-labels can be unreliable in terms of their large discrepancy from target ground truth. In this paper, we propose Cycle Self-Training (CST), a principled self-training algorithm that explicitly enforces pseudo-labels to generalize across domains. CST cycles between a forward step and a reverse step until convergence. In the forward step, CST generates target pseudo-labels with a source-trained classifier. In the reverse step, CST trains a target classifier using target pseudo-labels, and then updates the shared representations to make the target classifier perform well on the source data. We introduce the Tsallis entropy as a confidence-friendly regularization to improve the quality of target pseudo-labels. We analyze CST theoretically under realistic assumptions, and provide hard cases where CST recovers target ground truth, while both invariant feature learning and vanilla self-training fail. Empirical results indicate that CST significantly improves over the state-of-the-arts on visual recognition and sentiment analysis benchmarks.

[Novel Visual Category Discovery with Dual Ranking Statistics and Mutual Knowledge Distillation](#)

- Bingchen Zhao, Kai Han
- abstract@[open-review](#): In this paper, we tackle the problem of novel visual category discovery, i.e., grouping unlabelled images from new classes into different semantic partitions by leveraging a labelled dataset that contains images from other different but relevant categories. This is a more realistic and challenging setting than conventional semi-supervised learning. We propose a two-branch learning framework for this problem, with one branch focusing on local part-level information and the other branch focusing on overall characteristics. To transfer knowledge from the labelled data to the unlabelled, we propose using dual ranking statistics on both branches to generate pseudo labels for training on the unlabelled data. We further introduce a mutual knowledge distillation method to allow information exchange and encourage agreement between the two branches for discovering new categories, allowing our model to enjoy the benefits of global and local features. We comprehensively evaluate our method on public benchmarks for generic object classification, as well as the more challenging datasets for fine-grained visual recognition, achieving state-of-the-art performance.

[Stochastic Anderson Mixing for Nonconvex Stochastic Optimization](#)

- Fuchao Wei, Chenglong Bao, Yang Liu
- abstract@[open-review](#): Anderson mixing (AM) is an acceleration method for fixed-point iterations. Despite its success and wide usage in scientific computing, the convergence theory of AM remains unclear, and its applications to machine learning problems are not well explored. In this paper, by introducing damped projection and adaptive regularization to the classical AM, we propose a Stochastic Anderson Mixing (SAM) scheme to solve nonconvex stochastic optimization problems. Under mild assumptions, we establish the convergence theory of SAM, including the almost sure convergence to stationary points and the worst-case iteration complexity. Moreover, the complexity bound can be improved when randomly choosing an iterate as the output. To further accelerate the convergence, we incorporate a variance reduction technique into the proposed SAM. We also propose a preconditioned mixing strategy for SAM which can empirically achieve faster convergence or better generalization ability. Finally, we apply the SAM method to train various neural networks including the vanilla CNN, ResNets, WideResNet, ResNeXt, DenseNet and LSTM. Experimental results on image classification and language model demonstrate the advantages of our method.

[Sample-Efficient Reinforcement Learning for Linearly-Parameterized MDPs with a Generative Model](#)

- Bingyan Wang, Yuling Yan, Jianqing Fan
- abstract@[open-review](#): The curse of dimensionality is a widely known issue in reinforcement learning (RL). In the tabular setting where the state space \mathcal{S} and the action space \mathcal{A} are both finite, to obtain a near optimal policy with sampling access to a generative model, the minimax optimal sample complexity scales linearly with $|\mathcal{S}| \times |\mathcal{A}|$, which can be prohibitively large when $|\mathcal{S}|$ or $|\mathcal{A}|$ is large. This paper considers a Markov decision process (MDP) that admits a set of state-action features, which can linearly express (or approximate) its probability transition kernel. We show that a model-based approach (resp. $\sim Q$ -learning) provably learns an ε -optimal policy (resp. $\sim Q$ -function) with high probability as soon as the sample size exceeds the order of $\frac{K}{(1-\gamma)^3} \varepsilon^2$ (resp. $\frac{K}{(1-\gamma)^4} \varepsilon^2$), up to some logarithmic factor. Here K is the feature dimension and $\gamma \in (0, 1)$ is the discount factor of the MDP. Both sample complexity bounds are provably tight, and our result for the model-based approach matches the minimax lower bound. Our results show that for arbitrarily large-scale MDP, both the model-based approach and Q-learning are sample-efficient when K is relatively small, and hence the title of this paper.

[NN-Baker: A Neural-network Infused Algorithmic Framework for Optimization Problems on Geometric Intersection Graphs](#)

- Evan McCarty, Qi Zhao, Anastasios Sidiropoulos, Yusu Wang
- abstract@[open-review](#): Recent years have witnessed a surge of approaches to use neural networks to help tackle combinatorial optimization problems, including graph optimization problems. However, theoretical understanding of such approaches remains limited. In this paper, we consider the geometric setting, where graphs are induced by points in a fixed dimensional Euclidean space. We show that several graph optimization problems can be approximated by an algorithm that is polynomial in graph size n via a framework we propose, call the Baker-paradigm. More importantly, a key advantage of the Baker-paradigm is that it decomposes the input problem into (at most linear number of) small sub-problems of fixed sizes (independent of the size of the input). For the family of such fixed-size sub-problems, we can now design neural networks with universal approximation guarantees to solve them. This leads to a mixed algorithmic-ML framework, which we call NN-Baker that has the capacity to approximately solve a family of graph optimization problems (e.g, maximum independent set and minimum vertex cover) in time linear to input graph size, and only polynomial to approximation parameter. We instantiate our NN-Baker by a CNN version and GNN version, and demonstrate the effectiveness and efficiency of our approach via a range of experiments.

[A Note on Sparse Generalized Eigenvalue Problem](#)

- Yunfeng Cai, Guanhua Fang, Ping Li
- abstract@[open-review](#): The sparse generalized eigenvalue problem (SGEP) aims to find the leading eigenvector with sparsity structure. SGEP plays an important role in statistical learning and has wide applications including, but not limited to, sparse principal component analysis, sparse canonical correlation analysis and sparse Fisher discriminant analysis, etc. Due to the sparsity constraint, the solution of SGEP entails interesting properties from both numerical and statistical perspectives. In this paper, we provide a detailed sensitivity analysis for SGEP and establish the rate-optimal perturbation bound under the sparse setting. Specifically, we show that the bound is related to the perturbation/noise level and the recovery of the true support of the leading eigenvector as well. We also investigate the estimator of SGEP via imposing a non-convex regularization. Such estimator can achieve the optimal error rate and can recover the sparsity structure as well. Extensive numerical experiments corroborate our theoretical findings via using alternating direction method of multipliers (ADMM)-based computational method.

[RMIX: Learning Risk-Sensitive Policies for Cooperative Reinforcement Learning Agents](#)

- Wei Qiu, Xinrun Wang, Runsheng Yu, Rundong Wang, Xu He, Bo An, Svetlana Obraztsova, Zinovi Rabinovich
- abstract@[open-review](#): Current value-based multi-agent reinforcement learning methods optimize individual Q values to guide individuals' behaviours via centralized training with decentralized execution (CTDE). However, such expected, i.e., risk-neutral, Q value is not sufficient even with CTDE due to the

randomness of rewards and the uncertainty in environments, which causes the failure of these methods to train coordinating agents in complex environments. To address these issues, we propose RMIX, a novel cooperative MARL method with the Conditional Value at Risk (CVaR) measure over the learned distributions of individuals' Q values. Specifically, we first learn the return distributions of individuals to analytically calculate CVaR for decentralized execution. Then, to handle the temporal nature of the stochastic outcomes during executions, we propose a dynamic risk level predictor for risk level tuning. Finally, we optimize the CVaR policies with CVaR values used to estimate the target in TD error during centralized training and the CVaR values are used as auxiliary local rewards to update the local distribution via Quantile Regression loss. Empirically, we show that our method outperforms many state-of-the-art methods on various multi-agent risk-sensitive navigation scenarios and challenging StarCraft II cooperative tasks, demonstrating enhanced coordination and revealing improved sample efficiency.

Optimal Policies Tend To Seek Power

- Alex Turner, Logan Smith, Rohin Shah, Andrew Critch, Prasad Tadepalli
- abstract@[open-review](#): Some researchers speculate that intelligent reinforcement learning (RL) agents would be incentivized to seek resources and power in pursuit of the objectives we specify for them. Other researchers point out that RL agents need not have human-like power-seeking instincts. To clarify this discussion, we develop the first formal theory of the statistical tendencies of optimal policies. In the context of Markov decision processes, we prove that certain environmental symmetries are sufficient for optimal policies to tend to seek power over the environment. These symmetries exist in many environments in which the agent can be shut down or destroyed. We prove that in these environments, most reward functions make it optimal to seek power by keeping a range of options available and, when maximizing average reward, by navigating towards larger sets of potential terminal states.

Catalytic Role Of Noise And Necessity Of Inductive Biases In The Emergence Of Compositional Communication

- Łukasz Kuciński, Tomasz Korbak, Paweł Kołodziej, Piotr Miłoś
- abstract@[open-review](#): Communication is compositional if complex signals can be represented as a combination of simpler subparts. In this paper, we theoretically show that inductive biases on both the training framework and the data are needed to develop a compositional communication. Moreover, we prove that compositionality spontaneously arises in the signaling games, where agents communicate over a noisy channel. We experimentally confirm that a range of noise levels, which depends on the model and the data, indeed promotes compositionality. Finally, we provide a comprehensive study of this dependence and report results in terms of recently studied compositionality metrics: topographical similarity, conflict count, and context independence.

PLUR: A Unifying, Graph-Based View of Program Learning, Understanding, and Repair

- Zimin Chen, Vincent J Hellendoorn, Pascal Lamblin, Petros Maniatis, Pierre-Antoine Manzagol, Daniel Tarlow, Subhodeep Moitra
- abstract@[open-review](#): Machine learning for understanding and editing source code has recently attracted significant interest, with many developments in new models, new code representations, and new tasks. This proliferation can appear disparate and disconnected, making each approach seemingly unique and incompatible, thus obscuring the core machine learning challenges and contributions. In this work, we demonstrate that the landscape can be significantly simplified by taking a general approach of mapping a graph to a sequence of tokens and pointers. Our main result is to show that 16 recently published tasks of different shapes can be cast in this form, based on which a single model architecture achieves near or above state-of-the-art results on nearly all tasks, outperforming custom models like code2seq and alternative generic models like Transformers. This unification further enables multi-task learning and a series of cross-cutting experiments about the importance of different modeling choices for code understanding and repair tasks. The full framework, called PLUR, is easily extensible to more tasks, and will be open-sourced (<https://github.com/google-research/plur>).

COCO-LM: Correcting and Contrasting Text Sequences for Language Model Pretraining

- Yu Meng, Chenyan Xiong, Payal Bajaj, saurabh tiwary, Paul Bennett, Jiawei Han, XIA SONG
- abstract@[open-review](#): We present a self-supervised learning framework, COCO-LM, that pretrains Language Models by COrrecting and COntrasting corrupted text sequences. Following ELECTRA-style pretraining, COCO-LM employs an auxiliary language model to corrupt text sequences, upon which it constructs two new tasks for pretraining the main model. The first token-level task, Corrective Language Modeling, is to detect and correct tokens replaced by the auxiliary model, in order to better capture token-level semantics. The second sequence-level task, Sequence Contrastive Learning, is to align text sequences originated from the same source input while ensuring uniformity in the representation space. Experiments on GLUE and SQuAD demonstrate that COCO-LM not only outperforms recent state-of-the-art pretrained models in accuracy, but also improves pretraining efficiency. It achieves the MNLI accuracy of ELECTRA with 50% of its pretraining GPU hours. With the same pretraining steps of standard base/large-sized models, COCO-LM outperforms the previous best models by 1+ GLUE average points.

Minibatch and Momentum Model-based Methods for Stochastic Weakly Convex Optimization

- Qi Deng, Wenzhi Gao
- abstract@[open-review](#): Stochastic model-based methods have received increasing attention lately due to their appealing robustness to the stepsize selection and provable efficiency guarantee. We make two important extensions for improving model-based methods on stochastic weakly convex optimization. First, we propose new minibatch model-based methods by involving a set of samples to approximate the model function in each iteration. For the first time, we show that stochastic algorithms achieve linear speedup over the batch size even for non-smooth and non-convex (particularly, weakly convex) problems. To this end, we develop a novel sensitivity analysis of the proximal mapping involved in each algorithm iteration. Our analysis appears to be of independent interests in more general settings. Second, motivated by the success of momentum stochastic gradient descent, we propose a new stochastic extrapolated model-based method, greatly extending the classic Polyak momentum technique to a wider class of stochastic algorithms for weakly convex optimization. The rate of convergence to some natural stationarity condition is established over a fairly flexible range of extrapolation terms. While mainly focusing on weakly convex optimization, we also extend our work to convex optimization. We apply the minibatch and extrapolated model-based methods to stochastic convex optimization, for which we provide a new complexity bound and promising linear speedup in batch size. Moreover, an accelerated model-based method based on Nesterov's momentum is presented, for which we establish an optimal complexity bound for reaching optimality.

XDO: A Double Oracle Algorithm for Extensive-Form Games

- Stephen McAleer, JB Lanier, Kevin A Wang, Pierre Baldi, Roy Fox
- abstract@[open-review](#): Policy Space Response Oracles (PSRO) is a reinforcement learning (RL) algorithm for two-player zero-sum games that has been empirically shown to find approximate Nash equilibria in large games. Although PSRO is guaranteed to converge to an approximate Nash equilibrium and can handle continuous actions, it may take an exponential number of iterations as the number of information states (infostates) grows. We propose Extensive-Form Double Oracle (XDO), an extensive-form double oracle algorithm for two-player zero-sum games that is guaranteed to converge to an approximate Nash equilibrium linearly in the number of infostates. Unlike PSRO, which mixes best responses at the root of the game, XDO mixes best responses at every infostate. We also introduce Neural XDO (NXDO), where the best response is learned through deep RL. In tabular experiments on Leduc poker, we find that XDO achieves an approximate Nash equilibrium in a number of iterations an order of magnitude smaller than PSRO. Experiments on a modified Leduc poker game and Oshi-Zumo show that tabular XDO achieves a lower exploitability than CFR with the same amount of computation. We also find that NXDO outperforms PSRO and NFSP on a sequential multidimensional continuous-action game. NXDO is the first deep RL method that can find an approximate Nash equilibrium in high-dimensional continuous-action sequential games.

[Active Assessment of Prediction Services as Accuracy Surface Over Attribute Combinations](#)

- Vihari Piratla, Soumen Chakrabarti, Sunita Sarawagi
- abstract@[open-review](#): Our goal is to evaluate the accuracy of a black-box classification model, not as a single aggregate on a given test data distribution, but as a surface over a large number of combinations of attributes characterizing multiple test data distributions. Such attributed accuracy measures become important as machine learning models get deployed as a service, where the training data distribution is hidden from clients, and different clients may be interested in diverse regions of the data distribution. We present Attributed Accuracy Assay (AAA) --- a Gaussian Process (GP)-based probabilistic estimator for such an accuracy surface. Each attribute combination, called an 'arm' is associated with a Beta density from which the service's accuracy is sampled. We expect the GP to smooth the parameters of the Beta density over related arms to mitigate sparsity. We show that obvious application of GPs cannot address the challenge of heteroscedastic uncertainty over a huge attribute space that is sparsely and unevenly populated. In response, we present two enhancements: pooling sparse observations, and regularizing the scale parameter of the Beta densities. After introducing these innovations, we establish the effectiveness of AAA both in terms of its estimation accuracy and exploration efficiency, through extensive experiments and analysis.

[A mechanistic multi-area recurrent network model of decision-making](#)

- Michael Kleinman, Chandramouli Chandrasekaran, Jonathan Kao
- abstract@[open-review](#): Recurrent neural networks (RNNs) trained on neuroscience-based tasks have been widely used as models for cortical areas performing analogous tasks. However, very few tasks involve a single cortical area, and instead require the coordination of multiple brain areas. Despite the importance of multi-area computation, there is a limited understanding of the principles underlying such computation. We propose to use multi-area RNNs with neuroscience-inspired architecture constraints to derive key features of multi-area computation. In particular, we show that incorporating multiple areas and Dale's Law is critical for biasing the networks to learn biologically plausible solutions. Additionally, we leverage the full observability of the RNNs to show that output-relevant information is preferentially propagated between areas. These results suggest that cortex uses modular computation to generate minimal sufficient representations of task information. More broadly, our results suggest that constrained multi-area RNNs can produce experimentally testable hypotheses for computations that occur within and across multiple brain areas, enabling new insights into distributed computation in neural systems.

[Learning to Compose Visual Relations](#)

- Nan Liu, Shuang Li, Yilun Du, Josh Tenenbaum, Antonio Torralba
- abstract@[open-review](#): The visual world around us can be described as a structured set of objects and their associated relations. An image of a room may be conjured given only the description of the underlying objects and their associated relations. While there has been significant work on designing deep neural networks which may compose individual objects together, less work has been done on composing the individual relations between objects. A principal difficulty is that while the placement of objects is mutually independent, their relations are entangled and dependent on each other. To circumvent this issue, existing works primarily compose relations by utilizing a holistic encoder, in the form of text or graphs. In this work, we instead propose to represent each relation as an unnormalized density (an energy-based model), enabling us to compose separate relations in a factorized manner. We show that such a factorized decomposition allows the model to both generate and edit scenes that have multiple sets of relations more faithfully. We further show that decomposition enables our model to effectively understand the underlying relational scene structure.

[Identity testing for Mallows model](#)

- Róbert Busa-Fekete, Dimitris Fotakis, Balázs Szorenyi, Emmanouil Zampetakis
- abstract@[open-review](#): In this paper, we devise identity tests for ranking data that is generated from Mallows model both in the \emph{asymptotic} and \emph{non-asymptotic} settings. First we consider the case when the central ranking is known, and devise two algorithms for testing the spread parameter of the Mallows model. The first one is obtained by constructing a Uniformly Most Powerful Unbiased (UMPU) test in the asymptotic setting and then converting it into a sample-optimal non-asymptotic identity test. The resulting test is, however, impractical even for medium sized data, because it requires computing the distribution of the sufficient statistic. The second non-asymptotic test is derived from an optimal learning algorithm for the Mallows model. This test is both easy to compute and is sample-optimal for a wide range of parameters. Next, we consider testing Mallows models for the unknown central ranking case. This case can be tackled in the asymptotic setting by introducing a bias that exponentially decays with the sample size. We support all our findings with extensive numerical experiments and show that the proposed tests scale gracefully with the number of items to be ranked.

[Bandits with Knapsacks beyond the Worst Case](#)

- Karthik Abinav Sankararaman, Aleksandrs Slivkins
- abstract@[open-review](#): Bandits with Knapsacks (BwK) is a general model for multi-armed bandits under supply/budget constraints. While worst-case regret bounds for BwK are well-understood, we present three results that go beyond the worst-case perspective. First, we provide upper and lower bounds which amount to a full characterization for logarithmic, instance-dependent regret rates. Second, we consider "simple regret" in BwK, which tracks algorithm's performance in a given round, and prove that it is small in all but a few rounds. Third, we provide a "general reduction" from BwK to bandits which takes advantage of some known helpful structure, and apply this reduction to combinatorial semi-bandits, linear contextual bandits, and multinomial-logit bandits. Our results build on the BwK algorithm from prior work, providing new analyses thereof.

[Closing the loop in medical decision support by understanding clinical decision-making: A case study on organ transplantation](#)

- Yuchao Qin, Fergus Imrie, Alihan Hütük, Daniel Jarrett, Alexander Jimson, Mihaela van der Schaar
- abstract@[open-review](#): Significant effort has been placed on developing decision support tools to improve patient care. However, drivers of real-world clinical decisions in complex medical scenarios are not yet well-understood, resulting in substantial gaps between these tools and practical applications. In light of this, we highlight that more attention on understanding clinical decision-making is required both to elucidate current clinical practices and to enable effective human-machine interactions. This is imperative in high-stakes scenarios with scarce available resources. Using organ transplantation as a case study, we formalize the desiderata of methods for understanding clinical decision-making. We show that most existing machine learning methods are insufficient to meet these requirements and propose iTTransplant, a novel data-driven framework to learn the factors affecting decisions on organ offers in an instance-wise fashion directly from clinical data, as a possible solution. Through experiments on real-world liver transplantation data from OPTN, we demonstrate the use of iTTransplant to: (1) discover which criteria are most important to clinicians for organ offer acceptance; (2) identify patient-specific organ preferences of clinicians allowing automatic patient stratification; and (3) explore variations in transplantation practices between different transplant centers. Finally, we emphasize that the insights gained by iTTransplant can be used to inform the development of future decision support tools.

[Change Point Detection via Multivariate Singular Spectrum Analysis](#)

- Arwa Alanqary, Abdullah Alomar, Devavrat Shah
- abstract@[open-review](#): The objective of change point detection (CPD) is to detect significant and abrupt changes in the dynamics of the underlying system of interest through multivariate time series observations. In this work, we develop and analyze an algorithm for CPD that is inspired by a variant of the classical singular spectrum analysis (SSA) approach for time series by combining it with the classical cumulative sum (CUSUM) statistic from sequential

hypothesis testing. In particular, we model the underlying dynamics of multivariate time series observations through the spatio-temporal model introduced recently in the multivariate SSA (mSSA) literature. The change point in such a setting corresponds to a change in the underlying spatio-temporal model. As the primary contributions of this work, we develop an algorithm based on CUSUM-statistic to detect such change points in an online fashion. We extend the analysis of CUSUM statistics, traditionally done for the setting of independent observations, to the dependent setting of (multivariate) time series under the spatio-temporal model. Specifically, for a given parameter $\$h > 0\$$, our method achieves the following desirable trade-off: when a change happens, it detects it within $\$O(h)\$$ time delay on average, while in the absence of change, it does not declare false detection for at least $\$exp(\Omega(h))\$$ time length on average. We conduct empirical experiments using benchmark and synthetic datasets. We find that the proposed method performs competitively or outperforms the state-of-the-art change point detection methods across datasets.

Meta-learning to Improve Pre-training

- Aniruddh Raghu, Jonathan Lorraine, Simon Kornblith, Matthew McDermott, David K. Duvenaud
- abstract@[open-review](#): Pre-training (PT) followed by fine-tuning (FT) is an effective method for training neural networks, and has led to significant performance improvements in many domains. PT can incorporate various design choices such as task and data reweighting strategies, augmentation policies, and noise models, all of which can significantly impact the quality of representations learned. The hyperparameters introduced by these strategies therefore must be tuned appropriately. However, setting the values of these hyperparameters is challenging. Most existing methods either struggle to scale to high dimensions, are too slow and memory-intensive, or cannot be directly applied to the two-stage PT and FT learning process. In this work, we propose an efficient, gradient-based algorithm to meta-learn PT hyperparameters. We formalize the PT hyperparameter optimization problem and propose a novel method to obtain PT hyperparameter gradients by combining implicit differentiation and backpropagation through unrolled optimization. We demonstrate that our method improves predictive performance on two real-world domains. First, we optimize high-dimensional task weighting hyperparameters for multitask pre-training on protein-protein interaction graphs and improve AUROC by up to 3.9%. Second, we optimize a data augmentation neural network for self-supervised PT with SimCLR on electrocardiography data and improve AUROC by up to 1.9%.

Fair Sparse Regression with Clustering: An Invex Relaxation for a Combinatorial Problem

- Adarsh Barik, Jean Honorio
- abstract@[open-review](#): In this paper, we study the problem of fair sparse regression on a biased dataset where bias depends upon a hidden binary attribute. The presence of a hidden attribute adds an extra layer of complexity to the problem by combining sparse regression and clustering with unknown binary labels. The corresponding optimization problem is combinatorial, but we propose a novel relaxation of it as an invex optimization problem. To the best of our knowledge, this is the first invex relaxation for a combinatorial problem. We show that the inclusion of the debiasing/fairness constraint in our model has no adverse effect on the performance. Rather, it enables the recovery of the hidden attribute. The support of our recovered regression parameter vector matches exactly with the true parameter vector. Moreover, we simultaneously solve the clustering problem by recovering the exact value of the hidden attribute for each sample. Our method uses carefully constructed primal dual witnesses to provide theoretical guarantees for the combinatorial problem. To that end, we show that the sample complexity of our method is logarithmic in terms of the dimension of the regression parameter vector.

Probabilistic Margins for Instance Reweighting in Adversarial Training

- qizhou wang, Feng Liu, Bo Han, Tongliang Liu, Chen Gong, Gang Niu, Mingyuan Zhou, Masashi Sugiyama
- abstract@[open-review](#): Reweighting adversarial data during training has been recently shown to improve adversarial robustness, where data closer to the current decision boundaries are regarded as more critical and given larger weights. However, existing methods measuring the closeness are not very reliable: they are discrete and can take only a few values, and they are path-dependent, i.e., they may change given the same start and end points with different attack paths. In this paper, we propose three types of probabilistic margin (PM), which are continuous and path-independent, for measuring the aforementioned closeness and reweighing adversarial data. Specifically, a PM is defined as the difference between two estimated class-posterior probabilities, e.g., such a probability of the true label minus the probability of the most confusing label given some natural data. Though different PMs capture different geometric properties, all three PMs share a negative correlation with the vulnerability of data: data with larger/smaller PMs are safer/riskier and should have smaller/larger weights. Experiments demonstrated that PMs are reliable and PM-based reweighting methods outperformed state-of-the-art counterparts.

Unbalanced Optimal Transport through Non-negative Penalized Linear Regression

- Laetitia Chapel, Rémi Flamary, Haoran Wu, Cédric Févotte, Gilles Gasso
- abstract@[open-review](#): This paper addresses the problem of Unbalanced Optimal Transport (UOT) in which the marginal conditions are relaxed (using weighted penalties in lieu of equality) and no additional regularization is enforced on the OT plan. In this context, we show that the corresponding optimization problem can be reformulated as a non-negative penalized linear regression problem. This reformulation allows us to propose novel algorithms inspired from inverse problems and nonnegative matrix factorization. In particular, we consider majorization-minimization which leads in our setting to efficient multiplicative updates for a variety of penalties. Furthermore, we derive for the first time an efficient algorithm to compute the regularization path of UOT with quadratic penalties. The proposed algorithm provides a continuity of piece-wise linear OT plans converging to the solution of balanced OT (corresponding to infinite penalty weights). We perform several numerical experiments on simulated and real data illustrating the new algorithms, and provide a detailed discussion about more sophisticated optimization tools that can further be used to solve OT problems thanks to our reformulation.

The Difficulty of Passive Learning in Deep Reinforcement Learning

- Georg Ostrovski, Pablo Samuel Castro, Will Dabney
- abstract@[open-review](#): Learning to act from observational data without active environmental interaction is a well-known challenge in Reinforcement Learning (RL). Recent approaches involve constraints on the learned policy or conservative updates, preventing strong deviations from the state-action distribution of the dataset. Although these methods are evaluated using non-linear function approximation, theoretical justifications are mostly limited to the tabular or linear cases. Given the impressive results of deep reinforcement learning, we argue for a need to more clearly understand the challenges in this setting. In the vein of Held & Hein's classic 1963 experiment, we propose the "tandem learning" experimental paradigm which facilitates our empirical analysis of the difficulties in offline reinforcement learning. We identify function approximation in conjunction with fixed data distributions as the strongest factors, thereby extending but also challenging hypotheses stated in past work. Our results provide relevant insights for offline deep reinforcement learning, while also shedding new light on phenomena observed in the online case of learning control.

Intriguing Properties of Vision Transformers

- Muhammad Muzammal Naseer, Kanchana Ranasinghe, Salman H Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang
- abstract@[open-review](#): Vision transformers (ViT) have demonstrated impressive performance across numerous machine vision tasks. These models are based on multi-head self-attention mechanisms that can flexibly attend to a sequence of image patches to encode contextual cues. An important question is how such flexibility (in attending image-wide context conditioned on a given patch) can facilitate handling nuisances in natural images e.g., severe occlusions, domain shifts, spatial permutations, adversarial and natural perturbations. We systematically study this question via an extensive set of experiments encompassing three ViT families and provide comparisons with a high-performing convolutional neural network (CNN). We show and analyze the following intriguing properties of ViT: (a) Transformers are highly robust to severe occlusions, perturbations and domain shifts, e.g., retain as

high as 60% top-1 accuracy on ImageNet even after randomly occluding 80% of the image content. (b) The robustness towards occlusions is not due to texture bias, instead we show that ViTs are significantly less biased towards local textures, compared to CNNs. When properly trained to encode shape-based features, ViTs demonstrate shape recognition capability comparable to that of human visual system, previously unmatched in the literature. (c) Using ViTs to encode shape representation leads to an interesting consequence of accurate semantic segmentation without pixel-level supervision. (d) Off-the-shelf features from a single ViT model can be combined to create a feature ensemble, leading to high accuracy rates across a range of classification datasets in both traditional and few-shot learning paradigms. We show effective features of ViTs are due to flexible and dynamic receptive fields possible via self-attention mechanisms. Our code will be publicly released.

[PartialFed: Cross-Domain Personalized Federated Learning via Partial Initialization](#)

- Benyuan Sun, Hongxing Huo, YI YANG, Bo Bai
- abstract@[open-review](#): The burst of applications empowered by massive data have aroused unprecedented privacy concerns in AI society. Currently, data confidentiality protection has been one core issue during deep model training. Federated Learning (FL), which enables privacy-preserving training across multiple silos, gained rising popularity for its parameter-only communication. However, previous works have shown that FL revealed a significant performance drop if the data distributions are heterogeneous among different clients, especially when the clients have cross-domain characteristic, such as traffic, aerial and in-door. To address this challenging problem, we propose a novel idea, PartialFed, which loads a subset of the global model's parameters rather than loading the entire model used in most previous works. We first validate our algorithm with manually decided loading strategies inspired by various expert priors, named PartialFed-Fix. Then we develop PartialFed-Adaptive, which automatically selects personalized loading strategy for each client. The superiority of our algorithm is proved by demonstrating the new state-of-the-art results on cross-domain federated classification and detection. In particular, solely by initializing a small fraction of layers locally, we improve the performance of FedAvg on Office-Home and UODB by 4.88% and 2.65%, respectively. Further studies show that the adaptive strategy performs significantly better on domains with large deviation, e.g. improves AP50 by 4.03% and 4.89% on aerial and medical image detection compared to FedAvg.

[Adaptive Diffusion in Graph Neural Networks](#)

- Jialin Zhao, Yuxiao Dong, Ming Ding, Evgeny Kharlamov, Jie Tang
- abstract@[open-review](#): The success of graph neural networks (GNNs) largely relies on the process of aggregating information from neighbors defined by the input graph structures. Notably, message passing based GNNs, e.g., graph convolutional networks, leverage the immediate neighbors of each node during the aggregation process, and recently, graph diffusion convolution (GDC) is proposed to expand the propagation neighborhood by leveraging generalized graph diffusion. However, the neighborhood size in GDC is manually tuned for each graph by conducting grid search over the validation set, making its generalization practically limited. To address this issue, we propose the adaptive diffusion convolution (ADC) strategy to automatically learn the optimal neighborhood size from the data. Furthermore, we break the conventional assumption that all GNN layers and feature channels (dimensions) should use the same neighborhood for propagation. We design strategies to enable ADC to learn a dedicated propagation neighborhood for each GNN layer and each feature channel, making the GNN architecture fully coupled with graph structures--the unique property that differs GNNs from traditional neural networks. By directly plugging ADC into existing GNNs, we observe consistent and significant outperformance over both GDC and their vanilla versions across various datasets, demonstrating the improved model capacity brought by automatically learning unique neighborhood size per layer and per channel in GNNs.

[Recurrent Submodular Welfare and Matroid Blocking Semi-Bandits](#)

- Orestis Papadigenopoulos, Constantine Caramanis
- abstract@[open-review](#): A recent line of research focuses on the study of stochastic multi-armed bandits (MAB), in the case where temporal correlations of specific structure are imposed between the player's actions and the reward distributions of the arms. These correlations lead to (sub-)optimal solutions that exhibit interesting dynamical patterns -- a phenomenon that yields new challenges both from an algorithmic as well as a learning perspective. In this work, we extend the above direction to a combinatorial semi-bandit setting and study a variant of stochastic MAB, where arms are subject to matroid constraints and each arm becomes unavailable (blocked) for a fixed number of rounds after each play. A natural common generalization of the state-of-the-art for blocking bandits, and that for matroid bandits, only guarantees a $\frac{1}{2}$ -approximation for general matroids. In this paper we develop the novel technique of correlated (interleaved) scheduling, which allows us to obtain a polynomial-time $(1 - \frac{1}{e})$ -approximation algorithm (asymptotically and in expectation) for any matroid. Along the way, we discover an interesting connection to a variant of Submodular Welfare Maximization, for which we provide (asymptotically) matching upper and lower approximability bounds. In the case where the mean arm rewards are unknown, our technique naturally decouples the scheduling from the learning problem, and thus allows to control the $(1 - \frac{1}{e})$ -approximate regret of a UCB-based adaptation of our online algorithm.

[Representer Point Selection via Local Jacobian Expansion for Post-hoc Classifier Explanation of Deep Neural Networks and Ensemble Models](#)

- Yi Sui, Ga Wu, Scott Sanner
- abstract@[open-review](#): Explaining the influence of training data on deep neural network predictions is a critical tool for debugging models through data curation. A recent tractable and appealing approach for this task was provided via the concept of Representer Point Selection (RPS), i.e. a method that leverages the dual form of L_2 regularized optimization in the last layer of the neural network to identify the contribution of training points to the prediction. However, two key drawbacks of RPS are that they (i) lead to disagreement between the originally trained network and the RP regularized network modification and (ii) often yield a static ranking of training data for the same class, independent of the data being classified. Inspired by the RPS approach, we propose an alternative method based on a local Jacobian Taylor expansion (LJE) of the Jacobian. We empirically compared RPS-LJE with the original RPS- L_2 on image classification (with ResNet), text classification recurrent neural networks (with Bi-LSTM), and tabular classification (with XGBoost) tasks. Quantitatively, we show that RPS-LJE slightly outperforms RPS- L_2 and other state-of-the-art data explanation methods by up to 3% on a data debugging task. Qualitatively, we observe that RPS-LJE provides individualized explanations for each test data point rather than the class-specific static ranking of points in the original approach. Overall, RPS-LJE represents a novel approach to RPS that provides a powerful tool for data-oriented explanation and debugging.

[Editing a classifier by rewriting its prediction rules](#)

- Shibani Santurkar, Dimitris Tsipras, Mahalaxmi Elango, David Bau, Antonio Torralba, Aleksander Madry
- abstract@[open-review](#): We propose a methodology for modifying the behavior of a classifier by directly rewriting its prediction rules. Our method requires virtually no additional data collection and can be applied to a variety of settings, including adapting a model to new environments, and modifying it to ignore spurious features.

[How Modular should Neural Module Networks Be for Systematic Generalization?](#)

- Vanessa D'Amario, Tomotake Sasaki, Xavier Boix
- abstract@[open-review](#): Neural Module Networks (NMNs) aim at Visual Question Answering (VQA) via composition of modules that tackle a sub-task. NMNs are a promising strategy to achieve systematic generalization, i.e., overcoming biasing factors in the training distribution. However, the aspects of NMNs that facilitate systematic generalization are not fully understood. In this paper, we demonstrate that the degree of modularity of the NMN have

large influence on systematic generalization. In a series of experiments on three VQA datasets (VQA-MNIST, SQuOOP, and CLEVR-CoGenT), our results reveal that tuning the degree of modularity, especially at the image encoder stage, reaches substantially higher systematic generalization. These findings lead to new NMN architectures that outperform previous ones in terms of systematic generalization.

[Contrast and Mix: Temporal Contrastive Video Domain Adaptation with Background Mixing](#)

- Aadarsh Sahoo, Rutav Shah, Rameswar Panda, Kate Saenko, Abir Das
- abstract@[open-review](#): Unsupervised domain adaptation which aims to adapt models trained on a labeled source domain to a completely unlabeled target domain has attracted much attention in recent years. While many domain adaptation techniques have been proposed for images, the problem of unsupervised domain adaptation in videos remains largely underexplored. In this paper, we introduce Contrast and Mix (CoMix), a new contrastive learning framework that aims to learn discriminative invariant feature representations for unsupervised video domain adaptation. First, unlike existing methods that rely on adversarial learning for feature alignment, we utilize temporal contrastive learning to bridge the domain gap by maximizing the similarity between encoded representations of an unlabeled video at two different speeds as well as minimizing the similarity between different videos played at different speeds. Second, we propose a novel extension to the temporal contrastive loss by using background mixing that allows additional positives per anchor, thus adapting contrastive learning to leverage action semantics shared across both domains. Moreover, we also integrate a supervised contrastive learning objective using target pseudo-labels to enhance discriminability of the latent space for video domain adaptation. Extensive experiments on several benchmark datasets demonstrate the superiority of our proposed approach over state-of-the-art methods. Project page: <https://cvir.github.io/projects/comix>.

[The Flip Side of the Reweighted Coin: Duality of Adaptive Dropout and Regularization](#)

- Daniel LeJeune, Hamid Javadi, Richard Baraniuk
- abstract@[open-review](#): Among the most successful methods for sparsifying deep (neural) networks are those that adaptively mask the network weights throughout training. By examining this masking, or dropout, in the linear case, we uncover a duality between such adaptive methods and regularization through the so-called “ η -trick” that casts both as iteratively reweighted optimizations. We show that any dropout strategy that adapts to the weights in a monotonic way corresponds to an effective subquadratic regularization penalty, and therefore leads to sparse solutions. We obtain the effective penalties for several popular sparsification strategies, which are remarkably similar to classical penalties commonly used in sparse optimization. Considering variational dropout as a case study, we demonstrate similar empirical behavior between the adaptive dropout method and classical methods on the task of deep network sparsification, validating our theory.

[Active Learning of Convex Halfspaces on Graphs](#)

- Maximilian Thiessen, Thomas Gaertner
- abstract@[open-review](#): We systematically study the query complexity of learning geodesically convex halfspaces on graphs. Geodesic convexity is a natural generalisation of Euclidean convexity and allows the definition of convex sets and halfspaces on graphs. We prove an upper bound on the query complexity linear in the treewidth and the minimum hull set size but only logarithmic in the diameter. We show tight lower bounds along well-established separation axioms and identify the Radon number as a central parameter of the query complexity and the VC dimension. While previous bounds typically depend on the cut size of the labelling, all parameters in our bounds can be computed from the unlabelled graph. We provide evidence that ground-truth communities in real-world graphs are often convex and empirically compare our proposed approach with other active learning algorithms.

[Differentiable Spike: Rethinking Gradient-Descent for Training Spiking Neural Networks](#)

- Yuhang Li, Yufei Guo, Shanghang Zhang, Shikuang Deng, Yongqing Hai, Shi Gu
- abstract@[open-review](#): Spiking Neural Networks (SNNs) have emerged as a biology-inspired method mimicking the spiking nature of brain neurons. This bio-mimicry derives SNNs' energy efficiency of inference on neuromorphic hardware. However, it also causes an intrinsic disadvantage in training high-performing SNNs from scratch since the discrete spike prohibits the gradient calculation. To overcome this issue, the surrogate gradient (SG) approach has been proposed as a continuous relaxation. Yet the heuristic choice of SG leaves it vacant how the SG benefits the SNN training. In this work, we first theoretically study the gradient descent problem in SNN training and introduce finite difference gradient to quantitatively analyze the training behavior of SNN. Based on the introduced finite difference gradient, we propose a new family of Differentiable Spike (Dspike) functions that can adaptively evolve during training to find the optimal shape and smoothness for gradient estimation. Extensive experiments over several popular network structures show that training SNN with Dspike consistently outperforms the state-of-the-art training methods. For example, on the CIFAR10-DVS classification task, we can train a spiking ResNet-18 and achieve 75.4% top-1 accuracy with 10 time steps.

[Probabilistic Entity Representation Model for Reasoning over Knowledge Graphs](#)

- Nurendra Choudhary, Nikhil Rao, Sumeet Katariya, Karthik Subbian, Chandan Reddy
- abstract@[open-review](#): Logical reasoning over Knowledge Graphs (KGs) is a fundamental technique that can provide an efficient querying mechanism over large and incomplete databases. Current approaches employ spatial geometries such as boxes to learn query representations that encompass the answer entities and model the logical operations of projection and intersection. However, their geometry is restrictive and leads to non-smooth strict boundaries, which further results in ambiguous answer entities. Furthermore, previous works propose transformation tricks to handle unions which results in non-closure and, thus, cannot be chained in a stream. In this paper, we propose a Probabilistic Entity Representation Model (PERM) to encode entities as a Multivariate Gaussian density with mean and covariance parameters to capture its semantic position and smooth decision boundary, respectively. Additionally, we also define the closed logical operations of projection, intersection, and union that can be aggregated using an end-to-end objective function. On the logical query reasoning problem, we demonstrate that the proposed PERM significantly outperforms the state-of-the-art methods on various public benchmark KG datasets on standard evaluation metrics. We also evaluate PERM's competence on a COVID-19 drug-repurposing case study and show that our proposed work is able to recommend drugs with substantially better F1 than current methods. Finally, we demonstrate the working of our PERM's query answering process through a low-dimensional visualization of the Gaussian representations.

[Black Box Probabilistic Numerics](#)

- Onur Teymur, Christopher Foley, Philip Breen, Toni Karvonen, Chris J. Oates
- abstract@[open-review](#): Probabilistic numerics casts numerical tasks, such the numerical solution of differential equations, as inference problems to be solved. One approach is to model the unknown quantity of interest as a random variable, and to constrain this variable using data generated during the course of a traditional numerical method. However, data may be nonlinearly related to the quantity of interest, rendering the proper conditioning of random variables difficult and limiting the range of numerical tasks that can be addressed. Instead, this paper proposes to construct probabilistic numerical methods based only on the final output from a traditional method. A convergent sequence of approximations to the quantity of interest constitute a dataset, from which the limiting quantity of interest can be extrapolated, in a probabilistic analogue of Richardson's deferred approach to the limit. This black box approach (1) massively expands the range of tasks to which probabilistic numerics can be applied, (2) inherits the features and performance of state-of-the-art numerical methods, and (3) enables provably higher orders of convergence to be achieved. Applications are presented for nonlinear ordinary and partial differential equations, as well as for eigenvalue problems—a setting for which no probabilistic numerical methods have yet been developed.

Interpolation can hurt robust generalization even when there is no noise

- Konstantin Donhauser, Alexandru Tifrea, Michael Aerni, Reinhard Heckel, Fanny Yang
- abstract@[open-review](#): Numerous recent works show that overparameterization implicitly reduces variance for min-norm interpolators and max-margin classifiers. These findings suggest that ridge regularization has vanishing benefits in high dimensions. We challenge this narrative by showing that, even in the absence of noise, avoiding interpolation through ridge regularization can significantly improve generalization. We prove this phenomenon for the robust risk of both linear regression and classification, and hence provide the first theoretical result on \emph{robust overfitting}.

On the Equivalence between Neural Network and Support Vector Machine

- Yilan Chen, Wei Huang, Lam Nguyen, Tsui-Wei Weng
- abstract@[open-review](#): Recent research shows that the dynamics of an infinitely wide neural network (NN) trained by gradient descent can be characterized by Neural Tangent Kernel (NTK) \cite{jacot2018neural}. Under the squared loss, the infinite-width NN trained by gradient descent with an infinitely small learning rate is equivalent to kernel regression with NTK \cite{arora2019exact}. However, the equivalence is only known for ridge regression currently \cite{arora2019harnessing}, while the equivalence between NN and other kernel machines (KMs), e.g. support vector machine (SVM), remains unknown. Therefore, in this work, we propose to establish the equivalence between NN and SVM, and specifically, the infinitely wide NN trained by soft margin loss and the standard soft margin SVM with NTK trained by subgradient descent. Our main theoretical results include establishing the equivalence between NN and a broad family of ℓ_2 regularized KMs with finite-width bounds, which cannot be handled by prior work, and showing that every finite-width NN trained by such regularized loss functions is approximately a KM. Furthermore, we demonstrate our theory can enable three practical applications, including (i) \textit{non-vacuous} generalization bound of NN via the corresponding KM; (ii) \textit{nontrivial} robustness certificate for the infinite-width NN (while existing robustness verification methods would provide vacuous bounds); (iii) intrinsically more robust infinite-width NNs than those from previous kernel regression.

Learning Semantic Representations to Verify Hardware Designs

- Shobha Vasudevan, Wenjie (Joe) Jiang, David Bieber, Rishabh Singh, hamid shojaei, C. Richard Ho, Charles Sutton
- abstract@[open-review](#): Verification is a serious bottleneck in the industrial hardware design cycle, routinely requiring person-years of effort. Practical verification relies on a "best effort" process that simulates the design on test inputs. This suggests a new research question: Can this simulation data be exploited to learn a continuous representation of a hardware design that allows us to predict its functionality? As a first approach to this new problem, we introduce Design2Vec, a deep architecture that learns semantic abstractions of hardware designs. The key idea is to work at a higher level of abstraction than the gate or the bit level, namely the Register Transfer Level (RTL), which is somewhat analogous to software source code, and can be represented by a graph that incorporates control and data flow. This allows us to learn representations of RTL syntax and semantics using a graph neural network. We apply these representations to several tasks within verification, including predicting what cover points of the design will be exercised by a test, and generating new tests that will exercise desired cover points. We evaluate Design2Vec on three real-world hardware designs, including an industrial chip used in commercial data centers. Our results demonstrate that Design2Vec dramatically outperforms baseline approaches that do not incorporate the RTL semantics, scales to industrial designs, and can generate tests that exercise design points that are currently hard to cover with manually written tests by design verification experts.

Rebooting ACGAN: Auxiliary Classifier GANs with Stable Training

- Minguk Kang, Woohyeon Shim, Minsu Cho, Jaesik Park
- abstract@[open-review](#): Conditional Generative Adversarial Networks (cGAN) generate realistic images by incorporating class information into GAN. While one of the most popular cGANs is an auxiliary classifier GAN with softmax cross-entropy loss (ACGAN), it is widely known that training ACGAN is challenging as the number of classes in the dataset increases. ACGAN also tends to generate easily classifiable samples with a lack of diversity. In this paper, we introduce two cures for ACGAN. First, we identify that gradient exploding in the classifier can cause an undesirable collapse in early training, and projecting input vectors onto a unit hypersphere can resolve the problem. Second, we propose the Data-to-Data Cross-Entropy loss (D2D-CE) to exploit relational information in the class-labeled dataset. On this foundation, we propose the Rebooted Auxiliary Classifier Generative Adversarial Network (ReACGAN). The experimental results show that ReACGAN achieves state-of-the-art generation results on CIFAR10, Tiny-ImageNet, CUB200, and ImageNet datasets. We also verify that ReACGAN benefits from differentiable augmentations and that D2D-CE harmonizes with StyleGAN2 architecture. Model weights and a software package that provides implementations of representative cGANs and all experiments in our paper are available at <https://github.com/POSTECH-CVLab/PyTorch-StudioGAN>.

Towards a Theoretical Framework of Out-of-Distribution Generalization

- Haotian Ye, Chuanlong Xie, Tianle Cai, Ruichen Li, Zhenguo Li, Liwei Wang
- abstract@[open-review](#): Generalization to out-of-distribution (OOD) data is one of the central problems in modern machine learning. Recently, there is a surge of attempts to propose algorithms that mainly build upon the idea of extracting invariant features. Although intuitively reasonable, theoretical understanding of what kind of invariance can guarantee OOD generalization is still limited, and generalization to arbitrary out-of-distribution is clearly impossible. In this work, we take the first step towards rigorous and quantitative definitions of 1) what is OOD; and 2) what does it mean by saying an OOD problem is learnable. We also introduce a new concept of expansion function, which characterizes to what extent the variance is amplified in the test domains over the training domains, and therefore give a quantitative meaning of invariant features. Based on these, we prove an OOD generalization error bound. It turns out that OOD generalization largely depends on the expansion function. As recently pointed out by Gulrajani & Lopez-Paz (2020), any OOD learning algorithm without a model selection module is incomplete. Our theory naturally induces a model selection criterion. Extensive experiments on benchmark OOD datasets demonstrate that our model selection criterion has a significant advantage over baselines.

Slice Sampling Reparameterization Gradients

- David Zoltowski, Diana Cai, Ryan P. Adams
- abstract@[open-review](#): Many probabilistic modeling problems in machine learning use gradient-based optimization in which the objective takes the form of an expectation. These problems can be challenging when the parameters to be optimized determine the probability distribution under which the expectation is being taken, as the naive Monte Carlo procedure is not differentiable. Reparameterization gradients make it possible to efficiently perform optimization of these Monte Carlo objectives by transforming the expectation to be differentiable, but the approach is typically limited to distributions with simple forms and tractable normalization constants. Here we describe how to differentiate samples from slice sampling to compute \textit{slice sampling reparameterization gradients}, enabling a richer class of Monte Carlo objective functions to be optimized. Slice sampling is a Markov chain Monte Carlo algorithm for simulating samples from probability distributions; it only requires a density function that can be evaluated point-wise up to a normalization constant, making it applicable to a variety of inference problems and unnormalized models. Our approach is based on the observation that when the slice endpoints are known, the sampling path is a deterministic and differentiable function of the pseudo-random variables, since the algorithm is rejection-free. We evaluate the method on synthetic examples and apply it to a variety of applications with reparameterization of unnormalized probability distributions.

Multi-Label Learning with Pairwise Relevance Ordering

- Ming-Kun Xie, Sheng-Jun Huang
- abstract@[open-review](#): Precisely annotating objects with multiple labels is costly and has become a critical bottleneck in real-world multi-label classification tasks. Instead, deciding the relative order of label pairs is obviously less laborious than collecting exact labels. However, the supervised information of pairwise relevance ordering is less informative than exact labels. It is thus an important challenge to effectively learn with such weak supervision. In this paper, we formalize this problem as a novel learning framework, called multi-label learning with pairwise relevance ordering (PRO). We show that the unbiased estimator of classification risk can be derived with a cost-sensitive loss only from PRO examples. Theoretically, we provide the estimation error bound for the proposed estimator and further prove that it is consistent with respect to the commonly used ranking loss. Empirical studies on multiple datasets and metrics validate the effectiveness of the proposed method.

Sampling with Trustworthy Constraints: A Variational Gradient Framework

- Xingchao Liu, Xin Tong, Qiang Liu
- abstract@[open-review](#): Sampling-based inference and learning techniques, especially Bayesian inference, provide an essential approach to handling uncertainty in machine learning (ML). As these techniques are increasingly used in daily life, it becomes essential to safeguard the ML systems with various trustworthy-related constraints, such as fairness, safety, interpretability. Mathematically, enforcing these constraints in probabilistic inference can be cast into sampling from intractable distributions subject to general nonlinear constraints, for which practical efficient algorithms are still largely missing. In this work, we propose a family of constrained sampling algorithms which generalize Langevin Dynamics (LD) and Stein Variational Gradient Descent (SVGD) to incorporate a moment constraint specified by a general nonlinear function. By exploiting the gradient flow structure of LD and SVGD, we derive two types of algorithms for handling constraints, including a primal-dual gradient approach and the constraint controlled gradient descent approach. We investigate the continuous-time mean-field limit of these algorithms and show that they have $O(1/t)$ convergence under mild conditions. Moreover, the LD variant converges linearly assuming that a log Sobolev like inequality holds. Various numerical experiments are conducted to demonstrate the efficiency of our algorithms in trustworthy settings.

Robust and Decomposable Average Precision for Image Retrieval

- Elias Ramzi, Nicolas THOME, Clément Rambour, Nicolas Audebert, Xavier Bitot
- abstract@[open-review](#): In image retrieval, standard evaluation metrics rely on score ranking, e.g. average precision (AP). In this paper, we introduce a method for robust and decomposable average precision (ROADMAP) addressing two major challenges for end-to-end training of deep neural networks with AP: non-differentiability and non-decomposability. Firstly, we propose a new differentiable approximation of the rank function, which provides an upper bound of the AP loss and ensures robust training. Secondly, we design a simple yet effective loss function to reduce the decomposability gap between the AP in the whole training set and its averaged batch approximation, for which we provide theoretical guarantees. Extensive experiments conducted on three image retrieval datasets show that ROADMAP outperforms several recent AP approximation methods and highlight the importance of our two contributions. Finally, using ROADMAP for training deep models yields very good performances, outperforming state-of-the-art results on the three datasets. Code and instructions to reproduce our results will be made publicly available at <https://github.com/elias-ramzi/ROADMAP>.

Fast rates for prediction with limited expert advice

- El Mehdi Saad, Gilles Blanchard
- abstract@[open-review](#): We investigate the problem of minimizing the excess generalization error with respect to the best expert prediction in a finite family in the stochastic setting, under limited access to information. We consider that the learner has only access to a limited number of expert advices per training round, as well as for prediction. Assuming that the loss function is Lipschitz and strongly convex, we show that if we are allowed to see the advice of only one expert per round in the training phase, or to use the advice of only one expert for prediction in the test phase, the worst-case excess risk is $\$ \{ \Omega \} (1/\sqrt{T}) \$$ with probability lower bounded by a constant. However, if we are allowed to see at least two actively chosen expert advices per training round and use at least two experts for prediction, the fast rate $\$ \mathcal{O}(1/T) \$$ can be achieved. We design novel algorithms achieving this rate in this setting, and in the setting where the learner have a budget constraint on the total number of observed experts advices, and give precise instance-dependent bounds on the number of training rounds needed to achieve a given generalization error precision.

Probabilistic Transformer For Time Series Analysis

- Binh Tang, David S Matteson
- abstract@[open-review](#): Generative modeling of multivariate time series has remained challenging partly due to the complex, non-deterministic dynamics across long-distance timesteps. In this paper, we propose deep probabilistic methods that combine state-space models (SSMs) with transformer architectures. In contrast to previously proposed SSMs, our approaches use attention mechanism to model non-Markovian dynamics in the latent space and avoid recurrent neural networks entirely. We also extend our models to include several layers of stochastic variables organized in a hierarchy for further expressiveness. Compared to transformer models, ours are probabilistic, non-autoregressive, and capable of generating diverse long-term forecasts with uncertainty estimates. Extensive experiments show that our models consistently outperform competitive baselines on various tasks and datasets, including time series forecasting and human motion prediction.

A Hierarchical Reinforcement Learning Based Optimization Framework for Large-scale Dynamic Pickup and Delivery Problems

- Yi Ma, Xiaotian Hao, Jianye Hao, Jiawen Lu, Xing Liu, Tong Xialiang, Mingxuan Yuan, Zhigang Li, Jie Tang, Zhaopeng Meng
- abstract@[open-review](#): The Dynamic Pickup and Delivery Problem (DPDP) is an essential problem in the logistics domain, which is NP-hard. The objective is to dynamically schedule vehicles among multiple sites to serve the online generated orders such that the overall transportation cost could be minimized. The critical challenge of DPDP is the orders are not known a priori, i.e., the orders are dynamically generated in real-time. To address this problem, existing methods partition the overall DPDP into fixed-size sub-problems by caching online generated orders and solve each sub-problem, or on this basis to utilize the predicted future orders to optimize each sub-problem further. However, the solution quality and efficiency of these methods are unsatisfactory, especially when the problem scale is very large. In this paper, we propose a novel hierarchical optimization framework to better solve large-scale DPDPs. Specifically, we design an upper-level agent to dynamically partition the DPDP into a series of sub-problems with different scales to optimize vehicles routes towards globally better solutions. Besides, a lower-level agent is designed to efficiently solve each sub-problem by incorporating the strengths of classical operational research-based methods with reinforcement learning-based policies. To verify the effectiveness of the proposed framework, real historical data is collected from the order dispatching system of Huawei Supply Chain Business Unit and used to build a functional simulator. Extensive offline simulation and online testing conducted on the industrial order dispatching system justify the superior performance of our framework over existing baselines.

Spatio-Temporal Variational Gaussian Processes

- Oliver Hamelijnck, William Wilkinson, Niki Loppi, Arno Solin, Theodoros Damoulas
- abstract@[open-review](#): We introduce a scalable approach to Gaussian process inference that combines spatio-temporal filtering with natural gradient variational inference, resulting in a non-conjugate GP method for multivariate data that scales linearly with respect to time. Our natural gradient approach enables application of parallel filtering and smoothing, further reducing the temporal span complexity to be logarithmic in the number of time steps. We derive a sparse approximation that constructs a state-space model over a reduced set of spatial inducing points, and show that for separable Markov

kernels the full and sparse cases exactly recover the standard variational GP, whilst exhibiting favourable computational properties. To further improve the spatial scaling we propose a mean-field assumption of independence between spatial locations which, when coupled with sparsity and parallelisation, leads to an efficient and accurate method for large spatio-temporal problems.

[MERLOT: Multimodal Neural Script Knowledge Models](#)

- Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, Yejin Choi
- abstract@[open-review](#): As humans, we understand events in the visual world contextually, performing multimodal reasoning across time to make inferences about the past, present, and future. We introduce MERLOT, a model that learns multimodal script knowledge by watching millions of YouTube videos with transcribed speech -- in an entirely label-free, self-supervised manner. By pretraining with a mix of both frame-level (spatial) and video-level (temporal) objectives, our model not only learns to match images to temporally corresponding words, but also to contextualize what is happening globally over time. As a result, MERLOT exhibits strong out-of-the-box representations of temporal commonsense, and achieves state-of-the-art performance on 12 different video QA datasets when finetuned. It also transfers well to the world of static images, allowing models to reason about the dynamic context behind visual scenes. On Visual Commonsense Reasoning, MERLOT~answers questions correctly with 80.6% accuracy, outperforming state-of-the-art models of similar size by over 3%, even those that make heavy use of auxiliary supervised data (like object bounding boxes). Ablation analyses demonstrate the complementary importance of: 1) training on videos versus static images; 2) scaling the magnitude and diversity of the pretraining video corpus; and 3) using diverse objectives that encourage full-stack multimodal reasoning, from the recognition to cognition level.

[Fast Approximate Dynamic Programming for Infinite-Horizon Markov Decision Processes](#)

- Mohamad Amin Sharifi Kolarijani, Gyula Max, Peyman Mohajerin Mohajerin Esfahani
- abstract@[open-review](#): In this study, we consider the infinite-horizon, discounted cost, optimal control of stochastic nonlinear systems with separable cost and constraints in the state and input variables. Using the linear-time Legendre transform, we propose a novel numerical scheme for implementation of the corresponding value iteration (VI) algorithm in the conjugate domain. Detailed analyses of the convergence, time complexity, and error of the proposed algorithm are provided. In particular, with a discretization of size $\$X\$$ and $\$U\$$ for the state and input spaces, respectively, the proposed approach reduces the time complexity of each iteration in the VI algorithm from $\$O(XU)\$$ to $\$O(X+U)\$$, by replacing the minimization operation in the primal domain with a simple addition in the conjugate domain.

[Adaptive Risk Minimization: Learning to Adapt to Domain Shift](#)

- Marvin Zhang, Henrik Marklund, Nikita Dhawan, Abhishek Gupta, Sergey Levine, Chelsea Finn
- abstract@[open-review](#): A fundamental assumption of most machine learning algorithms is that the training and test data are drawn from the same underlying distribution. However, this assumption is violated in almost all practical applications: machine learning systems are regularly tested under distribution shift, due to changing temporal correlations, atypical end users, or other factors. In this work, we consider the problem setting of domain generalization, where the training data are structured into domains and there may be multiple test time shifts, corresponding to new domains or domain distributions. Most prior methods aim to learn a single robust model or invariant feature space that performs well on all domains. In contrast, we aim to learn models that adapt at test time to domain shift using unlabeled test points. Our primary contribution is to introduce the framework of adaptive risk minimization (ARM), in which models are directly optimized for effective adaptation to shift by learning to adapt on the training domains. Compared to prior methods for robustness, invariance, and adaptation, ARM methods provide performance gains of 1-4% test accuracy on a number of image classification problems exhibiting domain shift.

[Learning State Representations from Random Deep Action-conditional Predictions](#)

- Zeyu Zheng, Vivek Veeriah, Risto Vuorio, Richard L Lewis, Satinder Singh
- abstract@[open-review](#): Our main contribution in this work is an empirical finding that random General Value Functions (GVFs), i.e., deep action-conditional predictions--random both in what feature of observations they predict as well as in the sequence of actions the predictions are conditioned upon---form good auxiliary tasks for reinforcement learning (RL) problems. In particular, we show that random deep action-conditional predictions when used as auxiliary tasks yield state representations that produce control performance competitive with state-of-the-art hand-crafted auxiliary tasks like value prediction, pixel control, and CURL in both Atari and DeepMind Lab tasks. In another set of experiments we stop the gradients from the RL part of the network to the state representation learning part of the network and show, perhaps surprisingly, that the auxiliary tasks alone are sufficient to learn state representations good enough to outperform an end-to-end trained actor-critic baseline. We opensourced our code at https://github.com/Hwhitetooth/random_gvfs.

[Mixability made efficient: Fast online multiclass logistic regression](#)

- Rémi Jézéquel, Pierre Gaillard, Alessandro Rudi
- abstract@[open-review](#): Mixability has been shown to be a powerful tool to obtain algorithms with optimal regret. However, the resulting methods often suffer from high computational complexity which has reduced their practical applicability. For example, in the case of multiclass logistic regression, the aggregating forecaster (see Foster et al. 2018) achieves a regret of $\$O(\log(Bn))\$$ whereas Online Newton Step achieves $\$O(e^B \log(n))\$$ obtaining a double exponential gain in $\$B\$$ (a bound on the norm of comparative functions). However, this high statistical performance is at the price of a prohibitive computational complexity $\$O(n^{37})\$$. In this paper, we use quadratic surrogates to make aggregating forecasters more efficient. We show that the resulting algorithm has still high statistical performance for a large class of losses. In particular, we derive an algorithm for multiclass regression with a regret bounded by $\$O(B \log(n))\$$ and computational complexity of only $\$O(n^4)\$$.

[Tracking People with 3D Representations](#)

- Jathushan Rajasegaran, Georgios Pavlakos, Angjoo Kanazawa, Jitendra Malik
- abstract@[open-review](#): We present a novel approach for tracking multiple people in video. Unlike past approaches which employ 2D representations, we focus on using 3D representations of people, located in three-dimensional space. To this end, we develop a method, Human Mesh and Appearance Recovery (HMAR) which in addition to extracting the 3D geometry of the person as a SMPL mesh, also extracts appearance as a texture map on the triangles of the mesh. This serves as a 3D representation for appearance that is robust to viewpoint and pose changes. Given a video clip, we first detect bounding boxes corresponding to people, and for each one, we extract 3D appearance, pose, and location information using HMAR. These embedding vectors are then sent to a transformer, which performs spatio-temporal aggregation of the representations over the duration of the sequence. The similarity of the resulting representations is used to solve for associations that assigns each person to a tracklet. We evaluate our approach on the Posetrack, MuPoTs and AVA datasets. We find that 3D representations are more effective than 2D representations for tracking in these settings, and we obtain state-of-the-art performance. Code and results are available at: <https://brjathu.github.io/T3DP>.

[Off-Policy Risk Assessment in Contextual Bandits](#)

- Audrey Huang, Liu Leqi, Zachary Lipton, Kamyar Azizzadenesheli

- abstract@[open-review](#): Even when unable to run experiments, practitioners can evaluate prospective policies, using previously logged data. However, while the bandits literature has adopted a diverse set of objectives, most research on off-policy evaluation to date focuses on the expected reward. In this paper, we introduce Lipschitz risk functionals, a broad class of objectives that subsumes conditional value-at-risk (CVaR), variance, mean-variance, many distorted risks, and CPT risks, among others. We propose Off-Policy Risk Assessment (OPRA), a framework that first estimates a target policy's CDF and then generates plugin estimates for any collection of Lipschitz risks, providing finite sample guarantees that hold simultaneously over the entire class. We instantiate OPRA with both importance sampling and doubly robust estimators. Our primary theoretical contributions are (i) the first uniform concentration inequalities for both CDF estimators in contextual bandits and (ii) error bounds on our Lipschitz risk estimates, which all converge at a rate of $\$O(1/\sqrt{n})$.

[Adaptive Denoising via GainTuning](#)

- Sreyas Mohan, Joshua L Vincent, Ramon Manzorro, Peter Crozier, Carlos Fernandez-Granda, Eero Simoncelli
- abstract@[open-review](#): Deep convolutional neural networks (CNNs) for image denoising are typically trained on large datasets. These models achieve the current state of the art, but they do not generalize well to data that deviate from the training distribution. Recent work has shown that it is possible to train denoisers on a single noisy image. These models adapt to the features of the test image, but their performance is limited by the small amount of information used to train them. Here we propose "GainTuning", a methodology by which CNN models pre-trained on large datasets can be adaptively and selectively adjusted for individual test images. To avoid overfitting, GainTuning optimizes a single multiplicative scaling parameter (the "Gain") of each channel in the convolutional layers of the CNN. We show that GainTuning improves state-of-the-art CNNs on standard image-denoising benchmarks, boosting their denoising performance on nearly every image in a held-out test set. These adaptive improvements are even more substantial for test images differing systematically from the training data, either in noise level or image type. We illustrate the potential of adaptive GainTuning in a scientific application to transmission-electron-microscope images, using a CNN that is pre-trained on synthetic data. In contrast to the existing methodology, GainTuning is able to faithfully reconstruct the structure of catalytic nanoparticles from these data at extremely low signal-to-noise ratios.

[Optimal Sketching for Trace Estimation](#)

- Shuli Jiang, Hai Pham, David Woodruff, Richard Zhang
- abstract@[open-review](#): Matrix trace estimation is ubiquitous in machine learning applications and has traditionally relied on Hutchinson's method, which requires $\$O(\log(1/\delta)/\epsilon^2)$ matrix-vector product queries to achieve a $\$(1 \pm \epsilon)$ -multiplicative approximation to $\text{trace}(A)$ with failure probability δ on positive-semidefinite input matrices A . Recently, the Hutch++ algorithm was proposed, which reduces the number of matrix-vector queries from $\$O(1/\epsilon^2)$ to the optimal $\$O(1/\epsilon)$, and the algorithm succeeds with constant probability. However, in the high probability setting, the non-adaptive Hutch++ algorithm suffers an extra $\$O(\sqrt{\log(1/\delta)})$ multiplicative factor in its query complexity. Non-adaptive methods are important, as they correspond to sketching algorithms, which are mergeable, highly parallelizable, and provide low-memory streaming algorithms as well as low-communication distributed protocols. In this work, we close the gap between non-adaptive and adaptive algorithms, showing that even non-adaptive algorithms can achieve $\$O(\sqrt{\log(1/\delta)}\epsilon + \log(1/\delta))$ matrix-vector products. In addition, we prove matching lower bounds demonstrating that, up to a $\log \log(1/\delta)$ factor, no further improvement in the dependence on δ or ϵ is possible by any non-adaptive algorithm. Finally, our experiments demonstrate the superior performance of our sketch over the adaptive Hutch++ algorithm, which is less parallelizable, as well as over the non-adaptive Hutchinson's method.

[Estimating Multi-cause Treatment Effects via Single-cause Perturbation](#)

- Zhaozhi Qian, Alicia Curth, Mihaela van der Schaar
- abstract@[open-review](#): Most existing methods for conditional average treatment effect estimation are designed to estimate the effect of a single cause - only one variable can be intervened on at one time. However, many applications involve simultaneous intervention on multiple variables, which leads to multi-cause treatment effect problems. The multi-cause problem is challenging because one needs to overcome the confounding bias for a large number of treatment groups, each with a different cause combination. The combinatorial nature of the problem also leads to severe data scarcity - we only observe one factual outcome out of many potential outcomes. In this work, we propose Single-cause Perturbation (SCP), a novel two-step procedure to estimate the multi-cause treatment effect. SCP starts by augmenting the observational dataset with the estimated potential outcomes under single-cause interventions. It then performs covariate adjustment on the augmented dataset to obtain the estimator. SCP is agnostic to the exact choice of algorithm in either step. We show formally that the procedure is valid under standard assumptions in causal inference. We demonstrate the performance gain of SCP on extensive synthetic and semi-synthetic experiments.

[Be Confident! Towards Trustworthy Graph Neural Networks via Confidence Calibration](#)

- Xiao Wang, Hongrui Liu, Chuan Shi, Cheng Yang
- abstract@[open-review](#): Despite Graph Neural Networks (GNNs) have achieved remarkable accuracy, whether the results are trustworthy is still unexplored. Previous studies suggest that many modern neural networks are over-confident on the predictions, however, surprisingly, we discover that GNNs are primarily in the opposite direction, i.e., GNNs are under-confident. Therefore, the confidence calibration for GNNs is highly desired. In this paper, we propose a novel trustworthy GNN model by designing a topology-aware post-hoc calibration function. Specifically, we first verify that the confidence distribution in a graph has homophily property, and this finding inspires us to design a calibration GNN model (CaGCN) to learn the calibration function. CaGCN is able to obtain a unique transformation from logits of GNNs to the calibrated confidence for each node, meanwhile, such transformation is able to preserve the order between classes, satisfying the accuracy-preserving property. Moreover, we apply the calibration GNN to self-training framework, showing that more trustworthy pseudo labels can be obtained with the calibrated confidence and further improve the performance. Extensive experiments demonstrate the effectiveness of our proposed model in terms of both calibration and accuracy.

[Learning Riemannian metric for disease progression modeling](#)

- Samuel Gruffaz, Pierre-Emmanuel Poulet, Etienne Maheux, Bruno Jedynak, Stanley DURRLEMAN
- abstract@[open-review](#): Linear mixed-effect models provide a natural baseline for estimating disease progression using longitudinal data. They provide interpretable models at the cost of modeling assumptions on the progression profiles and their variability across subjects. A significant improvement is to embed the data in a Riemannian manifold and learn patient-specific trajectories distributed around a central geodesic. A few interpretable parameters characterize subject trajectories at the cost of a prior choice of the metric, which determines the shape of the trajectories. We extend this approach by learning the metric from the data allowing more flexibility while keeping the interpretability. Specifically, we learn the metric as the push-forward of the Euclidean metric by a diffeomorphism. This diffeomorphism is estimated iteratively as the composition of radial basis functions belonging to a reproducible kernel Hilbert space. The metric update allows us to improve the forecasting of imaging and clinical biomarkers in the Alzheimer's Disease Neuroimaging Initiative (ADNI) cohort. Our results compare favorably to the 56 methods benchmarked in the TADPOLE challenge.

[Bias and variance of the Bayesian-mean decoder](#)

- Arthur Prat-Carrabin, Michael Woodford
- abstract@[open-review](#): Perception, in theoretical neuroscience, has been modeled as the encoding of external stimuli into internal signals, which are then decoded. The Bayesian mean is an important decoder, as it is optimal for purposes of both estimation and discrimination. We present widely-applicable approximations to the bias and to the variance of the Bayesian mean, obtained under the minimal and biologically-relevant assumption that the encoding

results from a series of independent, though not necessarily identically-distributed, signals. Simulations substantiate the accuracy of our approximations in the small-noise regime. The bias of the Bayesian mean comprises two components: one driven by the prior, and one driven by the precision of the encoding. If the encoding is 'efficient', the two components have opposite effects; their relative strengths are determined by the objective that the encoding optimizes. The experimental literature on perception reports both 'Bayesian' biases directed towards prior expectations, and opposite, 'anti-Bayesian' biases. We show that different tasks are indeed predicted to yield such contradictory biases, under a consistently-optimal encoding-decoding model. Moreover, we recover Wei and Stocker's "law of human perception", a relation between the bias of the Bayesian mean and the derivative of its variance, and show how the coefficient of proportionality in this law depends on the task at hand. Our results provide a parsimonious theory of optimal perception under constraints, in which encoding and decoding are adapted both to the prior and to the task faced by the observer.

[MIRACLE: Causally-Aware Imputation via Learning Missing Data Mechanisms](#)

- Trent Kyono, Yao Zhang, Alexis Bellot, Mihaela van der Schaar
- abstract@[open-review](#): Missing data is an important problem in machine learning practice. Starting from the premise that imputation methods should preserve the causal structure of the data, we develop a regularization scheme that encourages any baseline imputation method to be causally consistent with the underlying data generating mechanism. Our proposal is a causally-aware imputation algorithm (MIRACLE). MIRACLE iteratively refines the imputation of a baseline by simultaneously modeling the missingness generating mechanism, encouraging imputation to be consistent with the causal structure of the data. We conduct extensive experiments on synthetic and a variety of publicly available datasets to show that MIRACLE is able to consistently improve imputation over a variety of benchmark methods across all three missingness scenarios: at random, completely at random, and not at random.

[Efficient Training of Visual Transformers with Small Datasets](#)

- Yahui Liu, Enver Sangineto, Wei Bi, Nicu Sebe, Bruno Lepri, Marco Nadai
- abstract@[open-review](#): Visual Transformers (VTs) are emerging as an architectural paradigm alternative to Convolutional networks (CNNs). Differently from CNNs, VTs can capture global relations between image elements and they potentially have a larger representation capacity. However, the lack of the typical convolutional inductive bias makes these models more data hungry than common CNNs. In fact, some local properties of the visual domain which are embedded in the CNN architectural design, in VTs should be learned from samples. In this paper, we empirically analyse different VTs, comparing their robustness in a small training set regime, and we show that, despite having a comparable accuracy when trained on ImageNet, their performance on smaller datasets can be largely different. Moreover, we propose an auxiliary self-supervised task which can extract additional information from images with only a negligible computational overhead. This task encourages the VTs to learn spatial relations within an image and makes the VT training much more robust when training data is scarce. Our task is used jointly with the standard (supervised) training and it does not depend on specific architectural choices, thus it can be easily plugged in the existing VTs. Using an extensive evaluation with different VTs and datasets, we show that our method can improve (sometimes dramatically) the final accuracy of the VTs. Our code is available at: <https://github.com/yhlleo/VTs-Drloc>.

[Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction](#)

- Dominik Stöger, Mahdi Soltanolkotabi
- abstract@[open-review](#): Recently there has been significant theoretical progress on understanding the convergence and generalization of gradient-based methods on nonconvex losses with overparameterized models. Nevertheless, many aspects of optimization and generalization and in particular the critical role of small random initialization are not fully understood. In this paper, we take a step towards demystifying this role by proving that small random initialization followed by a few iterations of gradient descent behaves akin to popular spectral methods. We also show that this implicit spectral bias from small random initialization, which is provably more prominent for overparameterized models, also puts the gradient descent iterations on a particular trajectory towards solutions that are not only globally optimal but also generalize well. Concretely, we focus on the problem of reconstructing a low-rank matrix from a few measurements via a natural nonconvex formulation. In this setting, we show that the trajectory of the gradient descent iterations from small random initialization can be approximately decomposed into three phases: (I) a spectral or alignment phase where we show that the iterates have an implicit spectral bias akin to spectral initialization allowing us to show that at the end of this phase the column space of the iterates and the underlying low-rank matrix are sufficiently aligned, (II) a saddle avoidance/refinement phase where we show that the trajectory of the gradient iterates moves away from certain degenerate saddle points, and (III) a local refinement phase where we show that after avoiding the saddles the iterates converge quickly to the underlying low-rank matrix. Underlying our analysis are insights for the analysis of overparameterized nonconvex optimization schemes that may have implications for computational problems beyond low-rank reconstruction.

[Efficient Combination of Rematerialization and Offloading for Training DNNs](#)

- Olivier Beaumont, Lionel Eyraud-Dubois, Alena Shilova
- abstract@[open-review](#): Rematerialization and offloading are two well known strategies to save memory during the training phase of deep neural networks, allowing data scientists to consider larger models, batch sizes or higher resolution data. Rematerialization trades memory for computation time, whereas Offloading trades memory for data movements. As these two resources are independent, it is appealing to consider the simultaneous combination of both strategies to save even more memory. We precisely model the costs and constraints corresponding to Deep Learning frameworks such as PyTorch or Tensorflow, we propose optimal algorithms to find a valid sequence of memory-constrained operations and finally, we evaluate the performance of proposed algorithms on realistic networks and computation platforms. Our experiments show that the possibility to offload can remove one third of the overhead of rematerialization, and that together they can reduce the memory used for activations by a factor 4 to 6, with an overhead below 20%.

[Particle Cloud Generation with Message Passing Generative Adversarial Networks](#)

- Raghav Kansal, Javier Duarte, Hao Su, Breno Orzari, Thiago Tomei, Maurizio Pierini, Mary Touranakou, Jean-roch Vlimant, Dimitrios Gunopulos
- abstract@[open-review](#): In high energy physics (HEP), jets are collections of correlated particles produced ubiquitously in particle collisions such as those at the CERN Large Hadron Collider (LHC). Machine learning (ML)-based generative models, such as generative adversarial networks (GANs), have the potential to significantly accelerate LHC jet simulations. However, despite jets having a natural representation as a set of particles in momentum-space, a.k.a. a particle cloud, there exist no generative models applied to such a dataset. In this work, we introduce a new particle cloud dataset (JetNet), and apply to it existing point cloud GANs. Results are evaluated using (1) 1-Wasserstein distances between high- and low-level feature distributions, (2) a newly developed Fréchet ParticleNet Distance, and (3) the coverage and (4) minimum matching distance metrics. Existing GANs are found to be inadequate for physics applications, hence we develop a new message passing GAN (MPGAN), which outperforms existing point cloud GANs on virtually every metric and shows promise for use in HEP. We propose JetNet as a novel point-cloud-style dataset for the ML community to experiment with, and set MPGAN as a benchmark to improve upon for future generative models. Additionally, to facilitate research and improve accessibility and reproducibility in this area, we release the open-source JetNet Python package with interfaces for particle cloud datasets, implementations for evaluation and loss metrics, and more tools for ML in HEP development.

[CoFiNet: Reliable Coarse-to-fine Correspondences for Robust PointCloud Registration](#)

- Hao Yu, Fu Li, Mahdi Saleh, Benjamin Busam, Slobodan Ilic

- abstract@[open-review](#): We study the problem of extracting correspondences between a pair of point clouds for registration. For correspondence retrieval, existing works benefit from matching sparse keypoints detected from dense points but usually struggle to guarantee their repeatability. To address this issue, we present CoFiNet - Coarse-to-Fine Network which extracts hierarchical correspondences from coarse to fine without keypoint detection. On a coarse scale and guided by a weighting scheme, our model firstly learns to match down-sampled nodes whose vicinity points share more overlap, which significantly shrinks the search space of a consecutive stage. On a finer scale, node proposals are consecutively expanded to patches that consist of groups of points together with associated descriptors. Point correspondences are then refined from the overlap areas of corresponding patches, by a density-adaptive matching module capable to deal with varying point density. Extensive evaluation of CoFiNet on both indoor and outdoor standard benchmarks shows our superiority over existing methods. Especially on 3DLoMatch where point clouds share less overlap, CoFiNet significantly outperforms state-of-the-art approaches by at least 5% on Registration Recall, with at most two-third of their parameters.

[Partial success in closing the gap between human and machine vision](#)

- Robert Geirhos, Kantharaju Narayananappa, Benjamin Mitzkus, Tizian Thieringer, Matthias Bethge, Felix A. Wichmann, Wieland Brendel
- abstract@[open-review](#): A few years ago, the first CNN surpassed human performance on ImageNet. However, it soon became clear that machines lack robustness on more challenging test cases, a major obstacle towards deploying machines "in the wild" and towards obtaining better computational models of human visual perception. Here we ask: Are we making progress in closing the gap between human and machine vision? To answer this question, we tested human observers on a broad range of out-of-distribution (OOD) datasets, recording 85,120 psychophysical trials across 90 participants. We then investigated a range of promising machine learning developments that crucially deviate from standard supervised CNNs along three axes: objective function (self-supervised, adversarially trained, CLIP language-image training), architecture (e.g. vision transformers), and dataset size (ranging from 1M to 1B). Our findings are threefold. (1.) The longstanding distortion robustness gap between humans and CNNs is closing, with the best models now exceeding human feedforward performance on most of the investigated OOD datasets. (2.) There is still a substantial image-level consistency gap, meaning that humans make different errors than models. In contrast, most models systematically agree in their categorisation errors, even substantially different ones like contrastive self-supervised vs. standard supervised models. (3.) In many cases, human-to-model consistency improves when training dataset size is increased by one to three orders of magnitude. Our results give reason for cautious optimism: While there is still much room for improvement, the behavioural difference between human and machine vision is narrowing. In order to measure future progress, 17 OOD datasets with image-level human behavioural data and evaluation code are provided as a toolbox and benchmark at: <https://github.com/bethgelab/model-vs-human/>

[LLC: Accurate, Multi-purpose Learnt Low-dimensional Binary Codes](#)

- Aditya Kusupati, Matthew Wallingford, Vivek Ramanujan, Raghav Somani, Jae Sung Park, Krishna Pillutla, Prateek Jain, Sham Kakade, Ali Farhadi
- abstract@[open-review](#): Learning binary representations of instances and classes is a classical problem with several high potential applications. In modern settings, the compression of high-dimensional neural representations to low-dimensional binary codes is a challenging task and often require large bit-codes to be accurate. In this work, we propose a novel method for learning low-dimensional binary codes for instances as well as classes. Our method does not require any side-information, like annotated attributes or label metadata, and learns extremely low-dimensional binary codes (approx 20 bits for ImageNet-1K). The learnt codes are super-efficient while still ensuring nearly optimal classification accuracy for ResNet50 on ImageNet-1K. We demonstrate that the learnt codes capture intrinsically important features in the data, by discovering an intuitive taxonomy over classes. We further quantitatively measure the quality of our codes by applying it to the efficient image retrieval as well as out-of-distribution (OOD) detection problems. For ImageNet-100 retrieval problem, our learnt binary codes outperform HashNet using only 10 bits and also are as accurate as 10 dimensional real representations. Finally, our learnt binary codes can perform OOD detection, out-of-the-box, as accurately as a baseline that needs 3000 samples to tune its threshold, while we require none. Code is open-sourced at <https://github.com/RAIVNLab/LLC>.

[Analytic Insights into Structure and Rank of Neural Network Hessian Maps](#)

- Sidak Pal Singh, Gregor Bachmann, Thomas Hofmann
- abstract@[open-review](#): The Hessian of a neural network captures parameter interactions through second-order derivatives of the loss. It is a fundamental object of study, closely tied to various problems in deep learning, including model design, optimization, and generalization. Most prior work has been empirical, typically focusing on low-rank approximations and heuristics that are blind to the network structure. In contrast, we develop theoretical tools to analyze the range of the Hessian map, which provide us with a precise understanding of its rank deficiency and the structural reasons behind it. This yields exact formulas and tight upper bounds for the Hessian rank of deep linear networks --- allowing for an elegant interpretation in terms of rank deficiency. Moreover, we demonstrate that our bounds remain faithful as an estimate of the numerical Hessian rank, for a larger class of models such as rectified and hyperbolic tangent networks. Further, we also investigate the implications of model architecture (e.g. width, depth, bias) on the rank deficiency. Overall, our work provides novel insights into the source and extent of redundancy in overparameterized neural networks.

[Well-tuned Simple Nets Excel on Tabular Datasets](#)

- Arlind Kadra, Marius Lindauer, Frank Hutter, Josif Grabocka
- abstract@[open-review](#): Tabular datasets are the last "unconquered castle" for deep learning, with traditional ML methods like Gradient-Boosted Decision Trees still performing strongly even against recent specialized neural architectures. In this paper, we hypothesize that the key to boosting the performance of neural networks lies in rethinking the joint and simultaneous application of a large set of modern regularization techniques. As a result, we propose regularizing plain Multilayer Perceptron (MLP) networks by searching for the optimal combination/cocktail of 13 regularization techniques for each dataset using a joint optimization over the decision on which regularizers to apply and their subsidiary hyperparameters. We empirically assess the impact of these regularization cocktails for MLPs in a large-scale empirical study comprising 40 tabular datasets and demonstrate that (i) well-regularized plain MLPs significantly outperform recent state-of-the-art specialized neural network architectures, and (ii) they even outperform strong traditional ML methods, such as XGBoost.

[POODLE: Improving Few-shot Learning via Penalizing Out-of-Distribution Samples](#)

- Duong Le, Khoi Duc Nguyen, Khoi Nguyen, Quoc-Huy Tran, Rang Nguyen, Binh-Son Hua
- abstract@[open-review](#): In this work, we propose to use out-of-distribution samples, i.e., unlabeled samples coming from outside the target classes, to improve few-shot learning. Specifically, we exploit the easily available out-of-distribution samples to drive the classifier to avoid irrelevant features by maximizing the distance from prototypes to out-of-distribution samples while minimizing that of in-distribution samples (i.e., support, query data). Our approach is simple to implement, agnostic to feature extractors, lightweight without any additional cost for pre-training, and applicable to both inductive and transductive settings. Extensive experiments on various standard benchmarks demonstrate that the proposed method consistently improves the performance of pretrained networks with different architectures.

[Combinatorial Pure Exploration with Bottleneck Reward Function](#)

- Yihan Du, Yuko Kuroki, Wei Chen
- abstract@[open-review](#): In this paper, we study the Combinatorial Pure Exploration problem with the Bottleneck reward function (CPE-B) under the fixed-confidence (FC) and fixed-budget (FB) settings. In CPE-B, given a set of base arms and a collection of subsets of base arms (super arms) following a certain combinatorial constraint, a learner sequentially plays a base arm and observes its random reward, with the objective of finding the optimal super

arm with the maximum bottleneck value, defined as the minimum expected reward of the base arms contained in the super arm.CPE-B captures a variety of practical scenarios such as network routing in communication networks, and its unique challenges fall on how to utilize the bottleneck property to save samples and achieve the statistical optimality. None of the existing CPE studies (most of them assume linear rewards) can be adapted to solve such challenges, and thus we develop brand-new techniques to handle them.For the FC setting, we propose novel algorithms with optimal sample complexity for a broad family of instances and establish a matching lower bound to demonstrate the optimality (within a logarithmic factor).For the FB setting, we design an algorithm which achieves the state-of-the-art error probability guarantee and is the first to run efficiently on fixed-budget path instances, compared to existing CPE algorithms. Our experimental results on the top-\$k\$, path and matching instances validate the empirical superiority of the proposed algorithms over their baselines.

[Densely connected normalizing flows](#)

- Matej Grcić, Ivan Grubišić, Siniša Šegvić
- abstract@[open-review](#): Normalizing flows are bijective mappings between inputs and latent representations with a fully factorized distribution. They are very attractive due to exact likelihood evaluation and efficient sampling. However, their effective capacity is often insufficient since the bijectivity constraint limits the model width. We address this issue by incrementally padding intermediate representations with noise. We precondition the noise in accordance with previous invertible units, which we describe as cross-unit coupling. Our invertible glow-like modules increase the model expressivity by fusing a densely connected block with Nyström self-attention. We refer to our architecture as DenseFlow since both cross-unit and intra-module couplings rely on dense connectivity. Experiments show significant improvements due to the proposed contributions and reveal state-of-the-art density estimation under moderate computing budgets.

[Snowflake: Scaling GNNs to high-dimensional continuous control via parameter freezing](#)

- Charles Blake, Vitaly Kurin, Maximilian Igl, Shimon Whiteson
- abstract@[open-review](#): Recent research has shown that graph neural networks (GNNs) can learn policies for locomotion control that are as effective as a typical multi-layer perceptron (MLP), with superior transfer and multi-task performance. However, results have so far been limited to training on small agents, with the performance of GNNs deteriorating rapidly as the number of sensors and actuators grows. A key motivation for the use of GNNs in the supervised learning setting is their applicability to large graphs, but this benefit has not yet been realised for locomotion control. We show that poor scaling in GNNs is a result of increasingly unstable policy updates, caused by overfitting in parts of the network during training. To combat this, we introduce Snowflake, a GNN training method for high-dimensional continuous control that freezes parameters in selected parts of the network. Snowflake significantly boosts the performance of GNNs for locomotion control on large agents, now matching the performance of MLPs while offering superior transfer properties.

[Subgame solving without common knowledge](#)

- Brian Zhang, Tuomas Sandholm
- abstract@[open-review](#): In imperfect-information games, subgame solving is significantly more challenging than in perfect-information games, but in the last few years, such techniques have been developed. They were the key ingredient to the milestone of superhuman play in no-limit Texas hold'em poker. Current subgame-solving techniques analyze the entire common-knowledge closure of the player's current information set, that is, the smallest set of nodes within which it is common knowledge that the current node lies. While this is acceptable in games like poker where the common-knowledge closure is relatively small, many practical games have more complex information structure, which renders the common-knowledge closure impractically large to enumerate or even reasonably approximate. We introduce an approach that overcomes this obstacle, by instead working with only low-order knowledge. Our approach allows an agent, upon arriving at an info-set, to basically prune any node that is no longer reachable, thereby massively reducing the game tree size relative to the common-knowledge subgame. We prove that, as is, our approach can increase exploitability compared to the blueprint strategy. However, we develop three avenues by which safety can be guaranteed. First, safety is guaranteed if the results of subgame solves are incorporated back into the blueprint. Second, we provide a method where safety is achieved by limiting the info-sets at which subgame solving is performed. Third, we prove that our approach, when applied at every info-set reached during play, achieves a weaker notion of equilibrium, which we coin affine equilibrium, and which may be of independent interest. We show that affine equilibria cannot be exploited by any Nash strategy of the opponent, so an opponent who wishes to exploit must open herself to counter-exploitation. Even without the safety-guaranteeing additions, experiments on medium-sized games show that our approach always reduced exploitability in practical games even when applied at every info-set, and a depth-limited version of it led to---to our knowledge---the first strong AI for the challenge problem dark chess.

[Fair Algorithms for Multi-Agent Multi-Armed Bandits](#)

- Safwan Hossain, Evi Micha, Nisarg Shah
- abstract@[open-review](#): We propose a multi-agent variant of the classical multi-armed bandit problem, in which there are \$N\$ agents and \$K\$ arms, and pulling an arm generates a (possibly different) stochastic reward for each agent. Unlike the classical multi-armed bandit problem, the goal is not to learn the "best arm"; indeed, each agent may perceive a different arm to be the best for her personally. Instead, we seek to learn a fair distribution over the arms. Drawing on a long line of research in economics and computer science, we use the Nash social welfare as our notion of fairness. We design multi-agent variants of three classic multi-armed bandit algorithms and show that they achieve sublinear regret, which is now measured in terms of the lost Nash social welfare. We also extend a classical lower bound, establishing the optimality of one of our algorithms.

[VAST: Value Function Factorization with Variable Agent Sub-Teams](#)

- Thomy Phan, Fabian Ritz, Lenz Belzner, Philipp Altmann, Thomas Gabor, Claudia Linnhoff-Popien
- abstract@[open-review](#): Value function factorization (VFF) is a popular approach to cooperative multi-agent reinforcement learning in order to learn local value functions from global rewards. However, state-of-the-art VFF is limited to a handful of agents in most domains. We hypothesize that this is due to the flat factorization scheme, where the VFF operator becomes a performance bottleneck with an increasing number of agents. Therefore, we propose VFF with variable agent sub-teams (VAST). VAST approximates a factorization for sub-teams which can be defined in an arbitrary way and vary over time, e.g., to adapt to different situations. The sub-team values are then linearly decomposed for all sub-team members. Thus, VAST can learn on a more focused and compact input representation of the original VFF operator. We evaluate VAST in three multi-agent domains and show that VAST can significantly outperform state-of-the-art VFF, when the number of agents is sufficiently large.

[On the Stochastic Stability of Deep Markov Models](#)

- Jan Drgona, Sayak Mukherjee, Jiaxin Zhang, Frank Liu, Mahantesh Halappanavar
- abstract@[open-review](#): Deep Markov models (DMM) are generative models which are scalable and expressive generalization of Markov models for representation, learning, and inference problems. However, the fundamental stochastic stability guarantees of such models have not been thoroughly investigated. In this paper, we present a novel stability analysis method and provide sufficient conditions of DMM's stochastic stability. The proposed stability analysis is based on the contraction of probabilistic maps modeled by deep neural networks. We make connections between the spectral properties of neural network's weights and different types of used activation function on the stability and overall dynamic behavior of DMMs with Gaussian distributions. Based on the theory, we propose a few practical methods for designing constrained DMMs with guaranteed stability. We empirically substantiate our theoretical results via intuitive numerical experiments using the proposed stability constraints.

[Multiwavelet-based Operator Learning for Differential Equations](#)

- Gaurav Gupta, Xiongye Xiao, Paul Bogdan
- abstract@[open-review](#): The solution of a partial differential equation can be obtained by computing the inverse operator map between the input and the solution space. Towards this end, we introduce a $\text{multiwavelet-based neural operator learning scheme}$ that compresses the associated operator's kernel using fine-grained wavelets. By explicitly embedding the inverse multiwavelet filters, we learn the projection of the kernel onto fixed multiwavelet polynomial bases. The projected kernel is trained at multiple scales derived from using repeated computation of multiwavelet transform. This allows learning the complex dependencies at various scales and results in a resolution-independent scheme. Compare to the prior works, we exploit the fundamental properties of the operator's kernel which enable numerically efficient representation. We perform experiments on the Korteweg-de Vries (KdV) equation, Burgers' equation, Darcy Flow, and Navier-Stokes equation. Compared with the existing neural operator approaches, our model shows significantly higher accuracy and achieves state-of-the-art in a range of datasets. For the time-varying equations, the proposed method exhibits a $(2X-10X)$ improvement (0.0018 (0.0033) relative $L2$ error for Burgers' (KdV) equation). By learning the mappings between function spaces, the proposed method has the ability to find the solution of a high-resolution input after learning from lower-resolution data.

[Intermediate Layers Matter in Momentum Contrastive Self Supervised Learning](#)

- Aakash Kaku, Sahana Upadhyaya, Narges Razavian
- abstract@[open-review](#): We show that bringing intermediate layers' representations of two augmented versions of an image closer together in self-supervised learning helps to improve the momentum contrastive (MoCo) method. To this end, in addition to the contrastive loss, we minimize the mean squared error between the intermediate layer representations or make their cross-correlation matrix closer to an identity matrix. Both loss objectives either outperform standard MoCo, or achieve similar performances on three diverse medical imaging datasets: NIH-Chest Xrays, Breast Cancer Histopathology, and Diabetic Retinopathy. The gains of the improved MoCo are especially large in a low-labeled data regime (e.g. 1% labeled data) with an average gain of 5% across three datasets. We analyze the models trained using our novel approach via feature similarity analysis and layer-wise probing. Our analysis reveals that models trained via our approach have higher feature reuse compared to a standard MoCo and learn informative features earlier in the network. Finally, by comparing the output probability distribution of models fine-tuned on small versus large labeled data, we conclude that our proposed method of pre-training leads to lower Kolmogorov–Smirnov distance, as compared to a standard MoCo. This provides additional evidence that our proposed method learns more informative features in the pre-training phase which could be leveraged in a low-labeled data regime.

[An Efficient Pessimistic-Optimistic Algorithm for Stochastic Linear Bandits with General Constraints](#)

- Xin Liu, Bin Li, Pengyi Shi, Lei Ying
- abstract@[open-review](#): This paper considers stochastic linear bandits with general nonlinear constraints. The objective is to maximize the expected cumulative reward over horizon T subject to a set of constraints in each round $\tau \leq T$. We propose a pessimistic-optimistic algorithm for this problem, which is efficient in two aspects. First, the algorithm yields $\tilde{\mathcal{O}}(\left(\frac{K^{0.75}}{\delta} + d\right)\sqrt{\tau})$ (pseudo) regret in round $\tau \leq T$, where K is the number of constraints, d is the dimension of the reward feature space, and δ is a Slater's constant; and ϵ constraint violation in any round $\tau > \tau'$, where τ' is ϵ independent of horizon T . Second, the algorithm is computationally efficient. Our algorithm is based on the primal-dual approach in optimization and includes two components. The primal component is similar to unconstrained stochastic linear bandits (our algorithm uses the linear upper confidence bound algorithm (LinUCB)). The computational complexity of the dual component depends on the number of constraints, but is independent of the sizes of the contextual space, the action space, and the feature space. Thus, the computational complexity of our algorithm is similar to LinUCB for unconstrained stochastic linear bandits.

[Efficiently Learning One Hidden Layer ReLU Networks From Queries](#)

- Sitan Chen, Adam Klivans, Raghu Meka
- abstract@[open-review](#): While the problem of PAC learning neural networks from samples has received considerable attention in recent years, in certain settings like model extraction attacks, it is reasonable to imagine having more than just the ability to observe random labeled examples. Motivated by this, we consider the following problem: given ϵ -black-box query access to a neural network F , recover F up to some error. Formally, we show that if F is an arbitrary one hidden layer neural network with ReLU activations, there is an algorithm with query complexity and runtime polynomial in all parameters which outputs a network F' achieving low square loss relative to F with respect to the Gaussian measure. While a number of works in the security literature have proposed and empirically demonstrated the effectiveness of certain algorithms for this problem, ours is to the best of our knowledge the first provable guarantee in this vein.

[Learning Nonparametric Volterra Kernels with Gaussian Processes](#)

- Magnus Ross, Michael T Smith, Mauricio Álvarez
- abstract@[open-review](#): This paper introduces a method for the nonparametric Bayesian learning of nonlinear operators, through the use of the Volterra series with kernels represented using Gaussian processes (GPs), which we term the nonparametric Volterra kernels model (NVKM). When the input function to the operator is unobserved and has a GP prior, the NVKM constitutes a powerful method for both single and multiple output regression, and can be viewed as a nonlinear and nonparametric latent force model. When the input function is observed, the NVKM can be used to perform Bayesian system identification. We use recent advances in efficient sampling of explicit functions from GPs to map process realisations through the Volterra series without resorting to numerical integration, allowing scalability through doubly stochastic variational inference, and avoiding the need for Gaussian approximations of the output processes. We demonstrate the performance of the model for both multiple output regression and system identification using standard benchmarks.

[DiBS: Differentiable Bayesian Structure Learning](#)

- Lars Lorch, Jonas Rothfuss, Bernhard Schölkopf, Andreas Krause
- abstract@[open-review](#): Bayesian structure learning allows inferring Bayesian network structure from data while reasoning about the epistemic uncertainty---a key element towards enabling active causal discovery and designing interventions in real world systems. In this work, we propose a general, fully differentiable framework for Bayesian structure learning (DiBS) that operates in the continuous space of a latent probabilistic graph representation. Contrary to existing work, DiBS is agnostic to the form of the local conditional distributions and allows for joint posterior inference of both the graph structure and the conditional distribution parameters. This makes our formulation directly applicable to posterior inference of nonstandard Bayesian network models, e.g., with nonlinear dependencies encoded by neural networks. Using DiBS, we devise an efficient, general purpose variational inference method for approximating distributions over structural models. In evaluations on simulated and real-world data, our method significantly outperforms related approaches to joint posterior inference.

[Nonparametric estimation of continuous DPPs with kernel methods](#)

- Michaël Fanuel, Rémi Bardenet
- abstract@[open-review](#): Determinantal Point Process (DPPs) are statistical models for repulsive point patterns. Both sampling and inference are tractable for DPPs, a rare feature among models with negative dependence that explains their popularity in machine learning and spatial statistics. Parametric and

nonparametric inference methods have been proposed in the finite case, i.e. when the point patterns live in a finite ground set. In the continuous case, only parametric methods have been investigated, while nonparametric maximum likelihood for DPPs -- an optimization problem over trace-class operators -- has remained an open question. In this paper, we show that a restricted version of this maximum likelihood (MLE) problem falls within the scope of a recent representer theorem for nonnegative functions in an RKHS. This leads to a finite-dimensional problem, with strong statistical ties to the original MLE. Moreover, we propose, analyze, and demonstrate a fixed point algorithm to solve this finite-dimensional problem. Finally, we also provide a controlled estimate of the correlation kernel of the DPP, thus providing more interpretability.

[FINE Samples for Learning with Noisy Labels](#)

- Taehyeon Kim, Jongwoo Ko, sangwook Cho, JinHwan Choi, Se-Young Yun
- abstract@[open-review](#): Modern deep neural networks (DNNs) become frail when the datasets contain noisy (incorrect) class labels. Robust techniques in the presence of noisy labels can be categorized into two folds: developing noise-robust functions or using noise-cleansing methods by detecting the noisy data. Recently, noise-cleansing methods have been considered as the most competitive noisy-label learning algorithms. Despite their success, their noisy label detectors are often based on heuristics more than a theory, requiring a robust classifier to predict the noisy data with loss values. In this paper, we propose a novel detector for filtering label noise. Unlike most existing methods, we focus on each data's latent representation dynamics and measure the alignment between the latent distribution and each representation using the eigen decomposition of the data gram matrix. Our framework, coined as filtering noisy instances via their eigenvectors (FINE), provides a robust detector with derivative-free simple methods having theoretical guarantees. Under our framework, we propose three applications of the FINE: sample-selection approach, semi-supervised learning approach, and collaboration with noise-robust loss functions. Experimental results show that the proposed methods consistently outperform corresponding baselines for all three applications on various benchmark datasets.

[Residual2Vec: Debiasing graph embedding with random graphs](#)

- Sadamori Kojaku, Jisung Yoon, Isabel Constantino, Yong-Yeol Ahn
- abstract@[open-review](#): Graph embedding maps a graph into a convenient vector-space representation for graph analysis and machine learning applications. Many graph embedding methods hinge on a sampling of context nodes based on random walks. However, random walks can be a biased sampler due to the structural properties of graphs. Most notably, random walks are biased by the degree of each node, where a node is sampled proportionally to its degree. The implication of such biases has not been clear, particularly in the context of graph representation learning. Here, we investigate the impact of the random walks' bias on graph embedding and propose residual2vec, a general graph embedding method that can debias various structural biases in graphs by using random graphs. We demonstrate that this debiasing not only improves link prediction and clustering performance but also allows us to explicitly model salient structural properties in graph embedding.

[Benign Overfitting in Multiclass Classification: All Roads Lead to Interpolation](#)

- Ke Wang, Vidya Muthukumar, Christos Thrampoulidis
- abstract@[open-review](#): The growing literature on "benign overfitting" in overparameterized models has been mostly restricted to regression or binary classification settings; however, most success stories of modern machine learning have been recorded in multiclass settings. Motivated by this discrepancy, we study benign overfitting in multiclass linear classification. Specifically, we consider the following popular training algorithms on separable data: (i) empirical risk minimization (ERM) with cross-entropy loss, which converges to the multiclass support vector machine (SVM) solution; (ii) ERM with least-squares loss, which converges to the min-norm interpolating (MNI) solution; and, (iii) the one-vs-all SVM classifier. Our first key finding is that under a simple sufficient condition, all three algorithms lead to classifiers that interpolate the training data and have equal accuracy. When the data is generated from Gaussian mixtures or a multinomial logistic model, this condition holds under high enough effective overparameterization. Second, we derive novel error bounds on the accuracy of the MNI classifier, thereby showing that all three training algorithms lead to benign overfitting under sufficient overparameterization. Ultimately, our analysis shows that good generalization is possible for SVM solutions beyond the realm in which typical margin-based bounds apply.

[Instance-Dependent Bounds for Zeroth-order Lipschitz Optimization with Error Certificates](#)

- Francois Bachoc, Tom Cesari, Sébastien Gerchinovitz
- abstract@[open-review](#): We study the problem of zeroth-order (black-box) optimization of a Lipschitz function f defined on a compact subset \mathcal{X} of \mathbb{R}^d , with the additional constraint that algorithms must certify the accuracy of their recommendations. We characterize the optimal number of evaluations of any Lipschitz function f to find and certify an approximate maximizer of f at accuracy ϵ . Under a weak assumption on \mathcal{X} , this optimal sample complexity is shown to be nearly proportional to the integral $\int_{\mathcal{X}} \mathcal{N}(x)^d dx / (\max(f) - f(x))^d$. This result, which was only (and partially) known in dimension $d=1$, solves an open problem dating back to 1991. In terms of techniques, our upper bound relies on a packing bound by Bouffier et al. (2020) for the Piyavskii-Shubert algorithm that we link to the above integral. We also show that a certified version of the computationally tractable DDO algorithm matches these packing and integral bounds. Our instance-dependent lower bound differs from traditional worst-case lower bounds in the Lipschitz setting and relies on a local worst-case analysis that could likely prove useful for other learning tasks.

[Training Neural Networks with Fixed Sparse Masks](#)

- Yi-Lin Sung, Varun Nair, Colin A. Raffel
- abstract@[open-review](#): During typical gradient-based training of deep neural networks, all of the model's parameters are updated at each iteration. Recent work has shown that it is possible to update only a small subset of the model's parameters during training, which can alleviate storage and communication requirements. In this paper, we show that it is possible to induce a fixed sparse mask on the model's parameters that selects a subset to update over many iterations. Our method constructs the mask out of the k parameters with the largest Fisher information as a simple approximation as to which parameters are most important for the task at hand. In experiments on parameter-efficient transfer learning and distributed training, we show that our approach matches or exceeds the performance of other methods for training with sparse updates while being more efficient in terms of memory usage and communication costs. We release our code publicly to promote further applications of our approach.

[VATT: Transformers for Multimodal Self-Supervised Learning from Raw Video, Audio and Text](#)

- Hassan Akbari, Liangzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, Boqing Gong
- abstract@[open-review](#): We present a framework for learning multimodal representations from unlabeled data using convolution-free Transformer architectures. Specifically, our Video-Audio-Text Transformer (VATT) takes raw signals as inputs and extracts multimodal representations that are rich enough to benefit a variety of downstream tasks. We train VATT end-to-end from scratch using multimodal contrastive losses and evaluate its performance by the downstream tasks of video action recognition, audio event classification, image classification, and text-to-video retrieval. Furthermore, we study a modality-agnostic single-backbone Transformer by sharing weights among the three modalities. We show that the convolution-free VATT outperforms state-of-the-art ConvNet-based architectures in the downstream tasks. Especially, VATT's vision Transformer achieves the top-1 accuracy of 82.1% on Kinetics-400, 83.6% on Kinetics-600, 72.7% on Kinetics-700, and 41.1% on Moments in Time, new records while avoiding supervised pre-training. Transferring to image classification leads to 78.7% top-1 accuracy on ImageNet compared to 64.7% by training the same Transformer from scratch,

showing the generalizability of our model despite the domain gap between videos and images. VATT's audio Transformer also sets a new record on waveform-based audio event recognition by achieving the mAP of 39.4% on AudioSet without any supervised pre-training.

Analyzing the Generalization Capability of SGLD Using Properties of Gaussian Channels

- Hao Wang, Yizhe Huang, Rui Gao, Flavio Calmon
- abstract@[open-review](#): Optimization is a key component for training machine learning models and has a strong impact on their generalization. In this paper, we consider a particular optimization method---the stochastic gradient Langevin dynamics (SGLD) algorithm---and investigate the generalization of models trained by SGLD. We derive a new generalization bound by connecting SGLD with Gaussian channels found in information and communication theory. Our bound can be computed from the training data and incorporates the variance of gradients for quantifying a particular kind of "sharpness" of the loss landscape. We also consider a closely related algorithm with SGLD, namely differentially private SGD (DP-SGD). We prove that the generalization capability of DP-SGD can be amplified by iteration. Specifically, our bound can be sharpened by including a time-decaying factor if the DP-SGD algorithm outputs the last iterate while keeping other iterates hidden. This decay factor enables the contribution of early iterations to our bound to reduce with time and is established by strong data processing inequalities---a fundamental tool in information theory. We demonstrate our bound through numerical experiments, showing that it can predict the behavior of the true generalization gap.

Learning to Schedule Heuristics in Branch and Bound

- Antonia Chmiela, Elias Khalil, Ambros Gleixner, Andrea Lodi, Sebastian Pokutta
- abstract@[open-review](#): Primal heuristics play a crucial role in exact solvers for Mixed Integer Programming (MIP). While solvers are guaranteed to find optimal solutions given sufficient time, real-world applications typically require finding good solutions early on in the search to enable fast decision-making. While much of MIP research focuses on designing effective heuristics, the question of how to manage multiple MIP heuristics in a solver has not received equal attention. Generally, solvers follow hard-coded rules derived from empirical testing on broad sets of instances. Since the performance of heuristics is problem-dependent, using these general rules for a particular problem might not yield the best performance. In this work, we propose the first data-driven framework for scheduling heuristics in an exact MIP solver. By learning from data describing the performance of primal heuristics, we obtain a problem-specific schedule of heuristics that collectively find many solutions at minimal cost. We formalize the learning task and propose an efficient algorithm for computing such a schedule. Compared to the default settings of a state-of-the-art academic MIP solver, we are able to reduce the average primal integral by up to 49% on two classes of challenging instances.

On Training Implicit Models

- Zhengyang Geng, Xin-Yu Zhang, Shaojie Bai, Yisen Wang, Zhouchen Lin
- abstract@[open-review](#): This paper focuses on training implicit models of infinite layers. Specifically, previous works employ implicit differentiation and solve the exact gradient for the backward propagation. However, is it necessary to compute such an exact but expensive gradient for training? In this work, we propose a novel gradient estimate for implicit models, named phantom gradient, that 1) forgoes the costly computation of the exact gradient; and 2) provides an update direction empirically preferable to the implicit model training. We theoretically analyze the condition under which an ascent direction of the loss landscape could be found and provide two specific instantiations of the phantom gradient based on the damped unrolling and Neumann series. Experiments on large-scale tasks demonstrate that these lightweight phantom gradients significantly accelerate the backward passes in training implicit models by roughly 1.7 \times and even boost the performance over approaches based on the exact gradient on ImageNet.

MLP-Mixer: An all-MLP Architecture for Vision

- Ilya O. Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, Alexey Dosovitskiy
- abstract@[open-review](#): Convolutional Neural Networks (CNNs) are the go-to model for computer vision. Recently, attention-based networks, such as the Vision Transformer, have also become popular. In this paper we show that while convolutions and attention are both sufficient for good performance, neither of them are necessary. We present MLP-Mixer, an architecture based exclusively on multi-layer perceptrons (MLPs). MLP-Mixer contains two types of layers: one with MLPs applied independently to image patches (i.e. "mixing" the per-location features), and one with MLPs applied across patches (i.e. "mixing" spatial information). When trained on large datasets, or with modern regularization schemes, MLP-Mixer attains competitive scores on image classification benchmarks, with pre-training and inference cost comparable to state-of-the-art models. We hope that these results spark further research beyond the realms of well established CNNs and Transformers.

A Framework to Learn with Interpretation

- Jayneel Parekh, Pavlo Mozharovskyi, Florence d'Alché-Buc
- abstract@[open-review](#): To tackle interpretability in deep learning, we present a novel framework to jointly learn a predictive model and its associated interpretation model. The interpreter provides both local and global interpretability about the predictive model in terms of human-understandable high level attribute functions, with minimal loss of accuracy. This is achieved by a dedicated architecture and well chosen regularization penalties. We seek for a small-size dictionary of high level attribute functions that take as inputs the outputs of selected hidden layers and whose outputs feed a linear classifier. We impose strong conciseness on the activation of attributes with an entropy-based criterion while enforcing fidelity to both inputs and outputs of the predictive model. A detailed pipeline to visualize the learnt features is also developed. Moreover, besides generating interpretable models by design, our approach can be specialized to provide post-hoc interpretations for a pre-trained neural network. We validate our approach against several state-of-the-art methods on multiple datasets and show its efficacy on both kinds of tasks.

One Loss for All: Deep Hashing with a Single Cosine Similarity based Learning Objective

- Jiun Tian Hoe, Kam Woh Ng, Tianyu Zhang, Chee Seng Chan, Yi-Zhe Song, Tao Xiang
- abstract@[open-review](#): A deep hashing model typically has two main learning objectives: to make the learned binary hash codes discriminative and to minimize a quantization error. With further constraints such as bit balance and code orthogonality, it is not uncommon for existing models to employ a large number (>4) of losses. This leads to difficulties in model training and subsequently impedes their effectiveness. In this work, we propose a novel deep hashing model with only $\{ \text{a single learning objective} \}$. Specifically, we show that maximizing the cosine similarity between the continuous codes and their corresponding $\{ \text{binary orthogonal codes} \}$ can ensure both hash code discriminativeness and quantization error minimization. Further, with this learning objective, code balancing can be achieved by simply using a Batch Normalization (BN) layer and multi-label classification is also straightforward with label smoothing. The result is a one-loss deep hashing model that removes all the hassles of tuning the weights of various losses. Importantly, extensive experiments show that our model is highly effective, outperforming the state-of-the-art multi-loss hashing models on three large-scale instance retrieval benchmarks, often by significant margins.

Fast and accurate randomized algorithms for low-rank tensor decompositions

- Linjian Ma, Edgar Solomonik

- abstract@[open-review](#): Low-rank Tucker and CP tensor decompositions are powerful tools in data analytics. The widely used alternating least squares (ALS) method, which solves a sequence of over-determined least squares subproblems, is costly for large and sparse tensors. We propose a fast and accurate sketched ALS algorithm for Tucker decomposition, which solves a sequence of sketched rank-constrained linear least squares subproblems. Theoretical sketch size upper bounds are provided to achieve $\mathcal{O}(\epsilon)$ relative error for each subproblem with two sketching techniques, TensorSketch and leverage score sampling. Experimental results show that this new ALS algorithm, combined with a new initialization scheme based on the randomized range finder, yields decomposition accuracy comparable to the standard higher-order orthogonal iteration (HOOI) algorithm. The new algorithm achieves up to 22.0% relative decomposition residual improvement compared to the state-of-the-art sketched randomized algorithm for Tucker decomposition of various synthetic and real datasets. This Tucker-ALS algorithm is further used to accelerate CP decomposition, by using randomized Tucker compression followed by CP decomposition of the Tucker core tensor. Experimental results show that this algorithm not only converges faster, but also yields more accurate CP decompositions.

[Communication-efficient SGD: From Local SGD to One-Shot Averaging](#)

- Artin Spiridonoff, Alex Olshevsky, Yannis Paschalidis
- abstract@[open-review](#): We consider speeding up stochastic gradient descent (SGD) by parallelizing it across multiple workers. We assume the same data set is shared among N workers, who can take SGD steps and coordinate with a central server. While it is possible to obtain a linear reduction in the variance by averaging all the stochastic gradients at every step, this requires a lot of communication between the workers and the server, which can dramatically reduce the gains from parallelism. The Local SGD method, proposed and analyzed in the earlier literature, suggests machines should make many local steps between such communications. While the initial analysis of Local SGD showed it needs $\Omega(\sqrt{T})$ communications for T local gradient steps in order for the error to scale proportionately to $1/(NT)$, this has been successively improved in a string of papers, with the state of the art requiring $\Omega(N \log T)$ communications. In this paper, we suggest a Local SGD scheme that communicates less overall by communicating less frequently as the number of iterations grows. Our analysis shows that this can achieve an error that scales as $1/(NT)$ with a number of communications that is completely independent of T . In particular, we show that $\Omega(N)$ communications are sufficient. Empirical evidence suggests this bound is close to tight as we further show that \sqrt{N} or $N^{3/4}$ communications fail to achieve linear speed-up in simulations. Moreover, we show that under mild assumptions, the main of which is twice differentiability on any neighborhood of the optimal solution, one-shot averaging which only uses a single round of communication can also achieve the optimal convergence rate asymptotically.

[Memory Efficient Meta-Learning with Large Images](#)

- John Bronskill, Daniela Massiceti, Massimiliano Patacchiola, Katja Hofmann, Sebastian Nowozin, Richard Turner
- abstract@[open-review](#): Meta learning approaches to few-shot classification are computationally efficient at test time, requiring just a few optimization steps or single forward pass to learn a new task, but they remain highly memory-intensive to train. This limitation arises because a task's entire support set, which can contain up to 1000 images, must be processed before an optimization step can be taken. Harnessing the performance gains offered by large images thus requires either parallelizing the meta-learner across multiple GPUs, which may not be available, or trade-offs between task and image size when memory constraints apply. We improve on both options by proposing LITE, a general and memory efficient episodic training scheme that enables meta-training on large tasks composed of large images on a single GPU. We achieve this by observing that the gradients for a task can be decomposed into a sum of gradients over the task's training images. This enables us to perform a forward pass on a task's entire training set but realize significant memory savings by back-propagating only a random subset of these images which we show is an unbiased approximation of the full gradient. We use LITE to train meta-learners and demonstrate new state-of-the-art accuracy on the real-world ORBIT benchmark and 3 of the 4 parts of the challenging VTAB+MD benchmark relative to leading meta-learners. LITE also enables meta-learners to be competitive with transfer learning approaches but at a fraction of the test-time computational cost, thus serving as a counterpoint to the recent narrative that transfer learning is all you need for few-shot classification.

[On the Power of Differentiable Learning versus PAC and SQ Learning](#)

- Emmanuel Abbe, Pritish Kamath, Eran Malach, Colin Sandon, Nathan Srebro
- abstract@[open-review](#): We study the power of learning via mini-batch stochastic gradient descent (SGD) on the loss of a differentiable model or neural network, and ask what learning problems can be learnt using this paradigm. We show that SGD can always simulate learning with statistical queries (SQ), but its ability to go beyond that depends on the precision ρ of the gradients and the minibatch size b . With fine enough precision relative to minibatch size, namely when $b \rho$ is small enough, SGD can go beyond SQ learning and simulate any sample-based learning algorithm and thus its learning power is equivalent to that of PAC learning; this extends prior work that achieved this result for $b=1$. Moreover, with polynomially many bits of precision (i.e. when ρ is exponentially small), SGD can simulate PAC learning regardless of the batch size. On the other hand, when $b \rho^2$ is large enough, the power of SGD is equivalent to that of SQ learning.

[Can we globally optimize cross-validation loss? Quasiconvexity in ridge regression](#)

- Will Stephenson, Zachary Frangella, Madeleine Udell, Tamara Broderick
- abstract@[open-review](#): Models like LASSO and ridge regression are extensively used in practice due to their interpretability, ease of use, and strong theoretical guarantees. Cross-validation (CV) is widely used for hyperparameter tuning in these models, but do practical methods minimize the true out-of-sample loss? A recent line of research promises to show that the optimum of the CV loss matches the optimum of the out-of-sample loss (possibly after simple corrections). It remains to show how tractable it is to minimize the CV loss. In the present paper, we show that, in the case of ridge regression, the CV loss may fail to be quasiconvex and thus may have multiple local optima. We can guarantee that the CV loss is quasiconvex in at least one case: when the spectrum of the covariate matrix is nearly flat and the noise in the observed responses is not too high. More generally, we show that quasiconvexity status is independent of many properties of the observed data (response norm, covariate-matrix right singular vectors and singular-value scaling) and has a complex dependence on the few that remain. We empirically confirm our theory using simulated experiments.

[Adaptive Proximal Gradient Methods for Structured Neural Networks](#)

- Jihun Yun, Aurelie C. Lozano, Eunho Yang
- abstract@[open-review](#): We consider the training of structured neural networks where the regularizer can be non-smooth and possibly non-convex. While popular machine learning libraries have resorted to stochastic (adaptive) subgradient approaches, the use of proximal gradient methods in the stochastic setting has been little explored and warrants further study, in particular regarding the incorporation of adaptivity. Towards this goal, we present a general framework of stochastic proximal gradient descent methods that allows for arbitrary positive preconditioners and lower semi-continuous regularizers. We derive two important instances of our framework: (i) the first proximal version of Adam, one of the most popular adaptive SGD algorithm, and (ii) a revised version of ProxQuant for quantization-specific regularizers, which improves upon the original approach by incorporating the effect of preconditioners in the proximal mapping computations. We provide convergence guarantees for our framework and show that adaptive gradient methods can have faster convergence in terms of constant than vanilla SGD for sparse data. Lastly, we demonstrate the superiority of stochastic proximal methods compared to subgradient-based approaches via extensive experiments. Interestingly, our results indicate that the benefit of proximal approaches over subgradient counterparts is more pronounced for non-convex regularizers than for convex ones.

[Discovering and Achieving Goals via World Models](#)

- Russell Mendonca, Oleh Rybkin, Kostas Daniilidis, Danijar Hafner, Deepak Pathak

- abstract@[open-review](#): How can artificial agents learn to solve many diverse tasks in complex visual environments without any supervision? We decompose this question into two challenges: discovering new goals and learning to reliably achieve them. Our proposed agent, Latent Explorer Achiever (LEXA), addresses both challenges by learning a world model from image inputs and using it to train an explorer and an achiever policy via imagined rollouts. Unlike prior methods that explore by reaching previously visited states, the explorer plans to discover unseen surprising states through foresight, which are then used as diverse targets for the achiever to practice. After the unsupervised phase, LEXA solves tasks specified as goal images zero-shot without any additional learning. LEXA substantially outperforms previous approaches to unsupervised goal reaching, both on prior benchmarks and on a new challenging benchmark with 40 test tasks spanning across four robotic manipulation and locomotion domains. LEXA further achieves goals that require interacting with multiple objects in sequence. Project page: <https://orybkin.github.io/lexa/>

Understanding and Improving Early Stopping for Learning with Noisy Labels

- Yingbin Bai, Erkun Yang, Bo Han, Yanhua Yang, Jiatong Li, Yinian Mao, Gang Niu, Tongliang Liu
- abstract@[open-review](#): The memorization effect of deep neural network (DNN) plays a pivotal role in many state-of-the-art label-noise learning methods. To exploit this property, the early stopping trick, which stops the optimization at the early stage of training, is usually adopted. Current methods generally decide the early stopping point by considering a DNN as a whole. However, a DNN can be considered as a composition of a series of layers, and we find that the latter layers in a DNN are much more sensitive to label noise, while their former counterparts are quite robust. Therefore, selecting a stopping point for the whole network may make different DNN layers antagonistically affect each other, thus degrading the final performance. In this paper, we propose to separate a DNN into different parts and progressively train them to address this problem. Instead of the early stopping which trains a whole DNN all at once, we initially train former DNN layers by optimizing the DNN with a relatively large number of epochs. During training, we progressively train the latter DNN layers by using a smaller number of epochs with the preceding layers fixed to counteract the impact of noisy labels. We term the proposed method as progressive early stopping (PES). Despite its simplicity, compared with the traditional early stopping, PES can help to obtain more promising and stable results. Furthermore, by combining PES with existing approaches on noisy label training, we achieve state-of-the-art performance on image classification benchmarks. The code is made public at <https://github.com/tmllab/PES>.

Distributionally Robust Imitation Learning

- Mohammad Ali Bashiri, Brian Ziebart, Xinhua Zhang
- abstract@[open-review](#): We consider the imitation learning problem of learning a policy in a Markov Decision Process (MDP) setting where the reward function is not given, but demonstrations from experts are available. Although the goal of imitation learning is to learn a policy that produces behaviors nearly as good as the experts' for a desired task, assumptions of consistent optimality for demonstrated behaviors are often violated in practice. Finding a policy that is distributionally robust against noisy demonstrations based on an adversarial construction potentially solves this problem by avoiding optimistic generalizations of the demonstrated data. This paper studies Distributionally Robust Imitation Learning (DRoIL) and establishes a close connection between DRoIL and Maximum Entropy Inverse Reinforcement Learning. We show that DRoIL can be seen as a framework that maximizes a generalized concept of entropy. We develop a novel approach to transform the objective function into a convex optimization problem over a polynomial number of variables for a class of loss functions that are additive over state and action spaces. Our approach lets us optimize both stationary and non-stationary policies and, unlike prevalent previous methods, it does not require repeatedly solving an inner reinforcement learning problem. We experimentally show the significant benefits of DRoIL's new optimization method on synthetic data and a highway driving environment.

On the Power of Edge Independent Graph Models

- Sudhanshu Chanpuriya, Cameron Musco, Konstantinos Sotiropoulos, Charalampos Tsourakakis
- abstract@[open-review](#): Why do many modern neural-network-based graph generative models fail to reproduce typical real-world network characteristics, such as high triangle density? In this work we study the limitations of \$edge\backslash independent\backslash random\backslash graph\backslash models\$, in which each edge is added to the graph independently with some probability. Such models include both the classic Erdos-Renyi and stochastic block models, as well as modern generative models such as NetGAN, variational graph autoencoders, and CELL. We prove that subject to a \$bounded\backslash overlap\$ condition, which ensures that the model does not simply memorize a single graph, edge independent models are inherently limited in their ability to generate graphs with high triangle and other subgraph densities. Notably, such high densities are known to appear in real-world social networks and other graphs. We complement our negative results with a simple generative model that balances overlap and accuracy, performing comparably to more complex models in reconstructing many graph statistics.

Stochastic Online Linear Regression: the Forward Algorithm to Replace Ridge

- Reda Ouhamma, Odalric-Ambrym Maillard, Vianney Perchet
- abstract@[open-review](#): We consider the problem of online linear regression in the stochastic setting. We derive high probability regret bounds for online \$\text{ridge}\$ regression and the \$\text{forward}\$ algorithm. This enables us to compare online regression algorithms more accurately and eliminate assumptions of bounded observations and predictions. Our study advocates for the use of the forward algorithm in lieu of ridge due to its enhanced bounds and robustness to the regularization parameter. Moreover, we explain how to integrate it in algorithms involving linear function approximation to remove a boundedness assumption without deteriorating theoretical bounds. We showcase this modification in linear bandit settings where it yields improved regret bounds. Last, we provide numerical experiments to illustrate our results and endorse our intuitions.

Dr Jekyll & Mr Hyde: the strange case of off-policy policy updates

- Romain Laroche, Remi Tachet des Combes
- abstract@[open-review](#): The policy gradient theorem states that the policy should only be updated in states that are visited by the current policy, which leads to insufficient planning in the off-policy states, and thus to convergence to suboptimal policies. We tackle this planning issue by extending the policy gradient theory to policy updates with respect to any state density. Under these generalized policy updates, we show convergence to optimality under a necessary and sufficient condition on the updates' state densities, and thereby solve the aforementioned planning issue. We also prove asymptotic convergence rates that significantly improve those in the policy gradient literature. To implement the principles prescribed by our theory, we propose an agent, Dr Jekyll & Mr Hyde (J&H), with a double personality: Dr Jekyll purely exploits while Mr Hyde purely explores. J&H's independent policies allow to record two separate replay buffers: one on-policy (Dr Jekyll's) and one off-policy (Mr Hyde's), and therefore to update J&H's models with a mixture of on-policy and off-policy updates. More than an algorithm, J&H defines principles for actor-critic algorithms to satisfy the requirements we identify in our analysis. We extensively test on finite MDPs where J&H demonstrates a superior ability to recover from converging to a suboptimal policy without impairing its speed of convergence. We also implement a deep version of the algorithm and test it on a simple problem where it shows promising results.

Understanding Adaptive, Multiscale Temporal Integration In Deep Speech Recognition Systems

- Menoua Keshishian, Samuel Norman-Hagnere, Nima Mesgarani
- abstract@[open-review](#): Natural signals such as speech are hierarchically structured across many different timescales, spanning tens (e.g., phonemes) to hundreds (e.g., words) of milliseconds, each of which is highly variable and context-dependent. While deep neural networks (DNNs) excel at recognizing complex patterns from natural signals, relatively little is known about how DNNs flexibly integrate across multiple timescales. Here, we show how a recently developed method for studying temporal integration in biological neural systems – the temporal context invariance (TCI) paradigm – can be used to understand temporal integration in DNNs. The method is simple: we measure responses to a large number of stimulus segments presented in two

different contexts and estimate the smallest segment duration needed to achieve a context invariant response. We applied our method to understand how the popular DeepSpeech2 model learns to integrate across time in speech. We find that nearly all of the model units, even in recurrent layers, have a compact integration window within which stimuli substantially alter the response and outside of which stimuli have little effect. We show that training causes these integration windows to shrink at early layers and expand at higher layers, creating a hierarchy of integration windows across the network. Moreover, by measuring integration windows for time-stretched/compressed speech, we reveal a transition point, midway through the trained network, where integration windows become yoked to the duration of stimulus structures (e.g., phonemes or words) rather than absolute time. Similar phenomena were observed in a purely recurrent and purely convolutional network although structure-yoked integration was more prominent in the recurrent network. These findings suggest that deep speech recognition systems use a common motif to encode the hierarchical structure of speech: integrating across short, time-yoked windows at early layers and long, structure-yoked windows at later layers. Our method provides a straightforward and general-purpose toolkit for understanding temporal integration in black-box machine learning models.

[VidLanKD: Improving Language Understanding via Video-Distilled Knowledge Transfer](#)

- Zineng Tang, Jaemin Cho, Hao Tan, Mohit Bansal
- abstract@[open-review](#): Since visual perception can give rich information beyond text descriptions for world understanding, there has been increasing interest in leveraging visual grounding for language learning. Recently, vokenization (Tan and Bansal, 2020) has attracted attention by using the predictions of a text-to-image retrieval model as labels for language model supervision. Despite its success, the method suffers from approximation error of using finite image labels and the lack of vocabulary diversity of a small image-text dataset. To overcome these limitations, we present VidLanKD, a video-language knowledge distillation method for improving language understanding. We train a multi-modal teacher model on a video-text dataset, and then transfer its knowledge to a student language model with a text dataset. To avoid approximation error, we propose to use different knowledge distillation objectives. In addition, the use of a large-scale video-text dataset helps learn diverse and richer vocabularies. In our experiments, VidLanKD achieves consistent improvements over text-only language models and vokenization models, on several downstream language understanding tasks including GLUE, SQuAD, and SWAG. We also demonstrate the improved world knowledge, physical reasoning, and temporal reasoning capabilities of our model by evaluating on the GLUE-diagnostics, PIQA, and TRACIE datasets. Lastly, we present comprehensive ablation studies as well as visualizations of the learned text-to-video grounding results of our teacher and student language models.

[Detecting Individual Decision-Making Style: Exploring Behavioral Stylometry in Chess](#)

- Reid McIlroy-Young, Yu Wang, Siddhartha Sen, Jon Kleinberg, Ashton Anderson
- abstract@[open-review](#): The advent of machine learning models that surpass human decision-making ability in complex domains has initiated a movement towards building AI systems that interact with humans. Many building blocks are essential for this activity, with a central one being the algorithmic characterization of human behavior. While much of the existing work focuses on aggregate human behavior, an important long-range goal is to develop behavioral models that specialize to individual people and can differentiate among them. To formalize this process, we study the problem of behavioral stylometry, in which the task is to identify a decision-maker from their decisions alone. We present a transformer-based approach to behavioral stylometry in the context of chess, where one attempts to identify the player who played a set of games. Our method operates in a few-shot classification framework, and can correctly identify a player from among thousands of candidate players with 98% accuracy given only 100 labeled games. Even when trained on amateur play, our method generalises to out-of-distribution samples of Grandmaster players, despite the dramatic differences between amateur and world-class players. Finally, we consider more broadly what our resulting embeddings reveal about human style in chess, as well as the potential ethical implications of powerful methods for identifying individuals from behavioral data.

[Coupled Gradient Estimators for Discrete Latent Variables](#)

- Zhe Dong, Andriy Mnih, George Tucker
- abstract@[open-review](#): Training models with discrete latent variables is challenging due to the high variance of unbiased gradient estimators. While low-variance reparameterization gradients of a continuous relaxation can provide an effective solution, a continuous relaxation is not always available or tractable. Dong et al. (2020) and Yin et al. (2020) introduced a performant estimator that does not rely on continuous relaxations; however, it is limited to binary random variables. We introduce a novel derivation of their estimator based on importance sampling and statistical couplings, which we extend to the categorical setting. Motivated by the construction of a stick-breaking coupling, we introduce gradient estimators based on reparameterizing categorical variables as sequences of binary variables and Rao-Blackwellization. In systematic experiments, we show that our proposed categorical gradient estimators provide state-of-the-art performance, whereas even with additional Rao-Blackwellization previous estimators (Yin et al., 2019) underperform a simpler REINFORCE with a leave-one-out-baseline estimator (Kool et al., 2019).

[AutoGEL: An Automated Graph Neural Network with Explicit Link Information](#)

- Zhili Wang, Shimin DI, Lei Chen
- abstract@[open-review](#): Recently, Graph Neural Networks (GNNs) have gained popularity in a variety of real-world scenarios. Despite the great success, the architecture design of GNNs heavily relies on manual labor. Thus, automated graph neural network (AutoGNN) has attracted interest and attention from the research community, which makes significant performance improvements in recent years. However, existing AutoGNN works mainly adopt an implicit way to model and leverage the link information in the graphs, which is not well regularized to the link prediction task on graphs, and limits the performance of AutoGNN for other graph tasks. In this paper, we present a novel AutoGNN work that explicitly models the link information, abbreviated to AutoGEL. In such a way, AutoGEL can handle the link prediction task and improve the performance of AutoGNNs on the node classification and graph classification task. Moreover, AutoGEL proposes a novel search space containing various design dimensions at both intra-layer and inter-layer designs and adopts a more robust differentiable search algorithm to further improve efficiency and effectiveness. Experimental results on benchmark data sets demonstrate the superiority of AutoGEL on several tasks.

[RL for Latent MDPs: Regret Guarantees and a Lower Bound](#)

- Jeongyeol Kwon, Yonathan Efroni, Constantine Caramanis, Shie Mannor
- abstract@[open-review](#): In this work, we consider the regret minimization problem for reinforcement learning in latent Markov Decision Processes (LMDP). In an LMDP, an MDP is randomly drawn from a set of $\$M\$$ possible MDPs at the beginning of the interaction, but the identity of the chosen MDP is not revealed to the agent. We first show that a general instance of LMDPs requires at least $\$|\Omega((SA)^M)|\$$ episodes to even approximate the optimal policy. Then, we consider sufficient assumptions under which learning good policies requires polynomial number of episodes. We show that the key link is a notion of separation between the MDP system dynamics. With sufficient separation, we provide an efficient algorithm with local guarantee, {\\it i.e.,} providing a sublinear regret guarantee when we are given a good initialization. Finally, if we are given standard statistical sufficiency assumptions common in the Predictive State Representation (PSR) literature (e.g., \cite{boots2011online}) and a reachability assumption, we show that the need for initialization can be removed.

[Adaptive Sampling for Minimax Fair Classification](#)

- Shubhangshu Shekhar, Greg Fields, Mohammad Ghavamzadeh, Tara Javidi
- abstract@[open-review](#): Machine learning models trained on uncurated datasets can often end up adversely affecting inputs belonging to underrepresented groups. To address this issue, we consider the problem of adaptively constructing training sets which allow us to learn classifiers that are fair in a {\em

minimax} sense. We first propose an adaptive sampling algorithm based on the principle of \emph{optimism}, and derive theoretical bounds on its performance. We also propose heuristic extensions of this algorithm suitable for application to large scale, practical problems. Next, by deriving algorithm independent lower-bounds for a specific class of problems, we show that the performance achieved by our adaptive scheme cannot be improved in general. We then validate the benefits of adaptively constructing training sets via experiments on synthetic tasks with logistic regression classifiers, as well as on several real-world tasks using convolutional neural networks (CNNs).

[Structured in Space, Randomized in Time: Leveraging Dropout in RNNs for Efficient Training](#)

- Anup Sarma, Sonali Singh, Huaipan Jiang, Rui Zhang, Mahmut Kandemir, Chita Das
- abstract@[open-review](#): Recurrent Neural Networks (RNNs), more specifically their Long Short-Term Memory (LSTM) variants, have been widely used as a deep learning tool for tackling sequence-based learning tasks in text and speech. Training of such LSTM applications is computationally intensive due to the recurrent nature of hidden state computation that repeats for each time step. While sparsity in Deep Neural Nets has been widely seen as an opportunity for reducing computation time in both training and inference phases, the usage of non-ReLU activation in LSTM RNNs renders the opportunities for such dynamic sparsity associated with neuron activation and gradient values to be limited or non-existent. In this work, we identify dropout induced sparsity for LSTMs as a suitable mode of computation reduction. Dropout is a widely used regularization mechanism, which randomly drops computed neuron values during each iteration of training. We propose to structure dropout patterns, by dropping out the same set of physical neurons within a batch, resulting in column (row) level hidden state sparsity, which are well amenable to computation reduction at run-time in general-purpose SIMD hardware as well as systolic arrays. We provide a detailed analysis of how the dropout-induced sparsity propagates through the different stages of network training and how it can be leveraged in each stage. More importantly, our proposed approach works as a direct replacement for existing dropout-based application settings. We conduct our experiments for three representative NLP tasks: language modelling on the PTB dataset, OpenNMT based machine translation using the IWSLT De-En and En-Vi datasets, and named entity recognition sequence labelling using the CoNLL-2003 shared task. We demonstrate that our proposed approach can be used to translate dropout-based computation reduction into reduced training time, with improvement ranging from 1.23\$\times\$ to 1.64\$\times\$, without sacrificing the target metric.

[Variational Continual Bayesian Meta-Learning](#)

- Qiang Zhang, Jinyuan Fang, Zaiqiao Meng, Shangsong Liang, Emine Yilmaz
- abstract@[open-review](#): Conventional meta-learning considers a set of tasks from a stationary distribution. In contrast, this paper focuses on a more complex online setting, where tasks arrive sequentially and follow a non-stationary distribution. Accordingly, we propose a Variational Continual Bayesian Meta-Learning (VC-BML) algorithm. VC-BML maintains a Dynamic Gaussian Mixture Model for meta-parameters, with the number of component distributions determined by a Chinese Restaurant Process. Dynamic mixtures at the meta-parameter level increase the capability to adapt to diverse tasks due to a larger parameter space, alleviating the negative knowledge transfer problem. To infer posteriors of model parameters, compared to the previously used point estimation method, we develop a more robust posterior approximation method -- structured variational inference for the sake of avoiding forgetting knowledge. Experiments on tasks from non-stationary distributions show that VC-BML is superior in transferring knowledge among diverse tasks and alleviating catastrophic forgetting in an online setting.

[Recognizing Vector Graphics without Rasterization](#)

- XINYANG JIANG, LU LIU, Caihua Shan, Yifei Shen, Xuanyi Dong, Dongsheng Li
- abstract@[open-review](#): In this paper, we consider a different data format for images: vector graphics. In contrast to raster graphics which are widely used in image recognition, vector graphics can be scaled up or down into any resolution without aliasing or information loss, due to the analytic representation of the primitives in the document. Furthermore, vector graphics are able to give extra structural information on how low-level elements group together to form high level shapes or structures. These merits of graphic vectors have not been fully leveraged in existing methods. To explore this data format, we target on the fundamental recognition tasks: object localization and classification. We propose an efficient CNN-free pipeline that does not render the graphic into pixels (i.e. rasterization), and takes textual document of the vector graphics as input, called YOLaT (You Only Look at Text). YOLaT builds multi-graphs to model the structural and spatial information in vector graphics, and a dual-stream graph neural network is proposed to detect objects from the graph. Our experiments show that by directly operating on vector graphics, YOLaT outperforms raster-graphic based object detection baselines in terms of both average precision and efficiency. Code is available at <https://github.com/microsoft/YOLaT-VectorGraphicsRecognition>.

[On Episodes, Prototypical Networks, and Few-Shot Learning](#)

- Steinar Laenen, Luca Bertinetto
- abstract@[open-review](#): Episodic learning is a popular practice among researchers and practitioners interested in few-shot learning. It consists of organising training in a series of learning problems (or episodes), each divided into a small training and validation subset to mimic the circumstances encountered during evaluation. But is this always necessary? In this paper, we investigate the usefulness of episodic learning in methods which use nonparametric approaches, such as nearest neighbours, at the level of the episode. For these methods, we not only show how the constraints imposed by episodic learning are not necessary, but that they in fact lead to a data-inefficient way of exploiting training batches. We conduct a wide range of ablative experiments with Matching and Prototypical Networks, two of the most popular methods that use nonparametric approaches at the level of the episode. Their "non-episodic" counterparts are considerably simpler, have less hyperparameters, and improve their performance in multiple few-shot classification datasets.

[Pointwise Bounds for Distribution Estimation under Communication Constraints](#)

- Wei-Ning Chen, Peter Kairouz, Ayfer Ozgur
- abstract@[open-review](#): We consider the problem of estimating a \$d\$-dimensional discrete distribution from its samples observed under a \$b\$-bit communication constraint. In contrast to most previous results that largely focus on the global minimax error, we study the local behavior of the estimation error and provide \emph{pointwise} bounds that depend on the target distribution \$p\$. In particular, we show that the \$\ell_2\$ error decays with \$O(\left(\frac{\|p\|}{\sqrt{n}}\right)^{1/2})\$ when \$n\$ is sufficiently large, hence it is governed by the \emph{half-norm} of \$p\$ instead of the ambient dimension \$d\$. For the achievability result, we propose a two-round sequentially interactive estimation scheme that achieves this error rate uniformly over all \$p\$. Our scheme is based on a novel local refinement idea, where we first use a standard global minimax scheme to localize \$p\$ and then use the remaining samples to locally refine our estimate. We also develop a new local minimax lower bound with (almost) matching \$\ell_2\$ error, showing that any interactive scheme must admit a \$\Omega(\left(\frac{\|p\|}{\sqrt{n}}\right)^{(1+\delta)/2})\$ \$\ell_2\$ error for any \$\delta > 0\$ when \$n\$ is sufficiently large. The lower bound is derived by first finding the best parametric sub-model containing \$p\$, and then upper bounding the quantized Fisher information under this model. Our upper and lower bounds together indicate that the \$\mathsf{H}(p) = \log(\sqrt{n}\|p\|)\$ bits of communication is both sufficient and necessary to achieve the optimal (centralized) performance, where \$\mathsf{H}(p)\$ is the Rényi entropy of order \$2\$. Therefore, under the \$\ell_2\$ loss, the correct measure of the local communication complexity at \$p\$ is its Rényi entropy.

[CHIP: CHannel Independence-based Pruning for Compact Neural Networks](#)

- Yang Sui, Miao Yin, Yi Xie, Huy Phan, Saman Aliari Zonouz, Bo Yuan
- abstract@[open-review](#): Filter pruning has been widely used for neural network compression because of its enabled practical acceleration. To date, most of the existing filter pruning works explore the importance of filters via using intra-channel information. In this paper, starting from an inter-channel

perspective, we propose to perform efficient filter pruning using Channel Independence, a metric that measures the correlations among different feature maps. The less independent feature map is interpreted as containing less useful information/knowledge, and hence its corresponding filter can be pruned without affecting model capacity. We systematically investigate the quantification metric, measuring scheme and sensitiveness/reliability of channel independence in the context of filter pruning. Our evaluation results for different models on various datasets show the superior performance of our approach. Notably, on CIFAR-10 dataset our solution can bring 0.75% and 0.94% accuracy increase over baseline ResNet-56 and ResNet-110 models, respectively, and meanwhile the model size and FLOPs are reduced by 42.8% and 47.4% (for ResNet-56) and 48.3% and 52.1% (for ResNet-110), respectively. On ImageNet dataset, our approach can achieve 40.8% and 44.8% storage and computation reductions, respectively, with 0.15% accuracy increase over the baseline ResNet-50 model. The code is available at https://github.com/Eclipsess/CHIP_NeurIPS2021.

Federated Split Task-Agnostic Vision Transformer for COVID-19 CXR Diagnosis

- Sangjoon Park, Gwanghyun Kim, Jeongsol Kim, Boah Kim, Jong Chul Ye
- abstract@[open-review](#): Federated learning, which shares the weights of the neural network across clients, is gaining attention in the healthcare sector as it enables training on a large corpus of decentralized data while maintaining data privacy. For example, this enables neural network training for COVID-19 diagnosis on chest X-ray (CXR) images without collecting patient CXR data across multiple hospitals. Unfortunately, the exchange of the weights quickly consumes the network bandwidth if highly expressive network architecture is employed. So-called split learning partially solves this problem by dividing a neural network into a client and a server part, so that the client part of the network takes up less extensive computation resources and bandwidth. However, it is not clear how to find the optimal split without sacrificing the overall network performance. To amalgamate these methods and thereby maximize their distinct strengths, here we show that the Vision Transformer, a recently developed deep learning architecture with straightforward decomposable configuration, is ideally suitable for split learning without sacrificing performance. Even under the non-independent and identically distributed data distribution which emulates a real collaboration between hospitals using CXR datasets from multiple sources, the proposed framework was able to attain performance comparable to data-centralized training. In addition, the proposed framework along with heterogeneous multi-task clients also improves individual task performances including the diagnosis of COVID-19, eliminating the need for sharing large weights with innumerable parameters. Our results affirm the suitability of Transformer for collaborative learning in medical imaging and pave the way forward for future real-world implementations.

Active Offline Policy Selection

- Ksenia Konyushova, Yutian Chen, Thomas Paine, Caglar Gulcehre, Cosmin Paduraru, Daniel J. Mankowitz, Misha Denil, Nando de Freitas
- abstract@[open-review](#): This paper addresses the problem of policy selection in domains with abundant logged data, but with a restricted interaction budget. Solving this problem would enable safe evaluation and deployment of offline reinforcement learning policies in industry, robotics, and recommendation domains among others. Several off-policy evaluation (OPE) techniques have been proposed to assess the value of policies using only logged data. However, there is still a big gap between the evaluation by OPE and the full online evaluation in the real environment. Yet, large amounts of online interactions are often not possible in practice. To overcome this problem, we introduce active offline policy selection --- a novel sequential decision approach that combines logged data with online interaction to identify the best policy. This approach uses OPE estimates to warm start the online evaluation. Then, in order to utilize the limited environment interactions wisely we decide which policy to evaluate next based on a Bayesian optimization method with a kernel function that represents policy similarity. We use multiple benchmarks with a large number of candidate policies to show that the proposed approach improves upon state-of-the-art OPE estimates and pure online policy evaluation.

Unsupervised Representation Transfer for Small Networks: I Believe I Can Distill On-the-Fly

- Hee Min Choi, Hyo Kang, Dokwan Oh
- abstract@[open-review](#): A current remarkable improvement of unsupervised visual representation learning is based on heavy networks with large-batch training. While recent methods have greatly reduced the gap between supervised and unsupervised performance of deep models such as ResNet-50, this development has been relatively limited for small models. In this work, we propose a novel unsupervised learning framework for small networks that combines deep self-supervised representation learning and knowledge distillation within one-phase training. In particular, a teacher model is trained to produce consistent cluster assignments between different views of the same image. Simultaneously, a student model is encouraged to mimic the prediction of on-the-fly self-supervised teacher. For effective knowledge transfer, we adopt the idea of domain classifier so that student training is guided by discriminative features invariant to the representational space shift between teacher and student. We also introduce a network driven multi-view generation paradigm to capture rich feature information contained in the network itself. Extensive experiments show that our student models surpass state-of-the-art offline distilled networks even from stronger self-supervised teachers as well as top-performing self-supervised models. Notably, our ResNet-18, trained with ResNet-50 teacher, achieves 68.3% ImageNet Top-1 accuracy on frozen feature linear evaluation, which is only 1.5% below the supervised baseline.

Understanding Bandits with Graph Feedback

- Houshuang Chen, zengfeng Huang, Shuai Li, Chihao Zhang
- abstract@[open-review](#): The bandit problem with graph feedback, proposed in [Mannor and Shamir, NeurIPS 2011], is modeled by a directed graph $G = (V, E)$ where V is the collection of bandit arms, and once an arm is triggered, all its incident arms are observed. A fundamental question is how the structure of the graph affects the min-max regret. We propose the notions of the fractional weak domination number δ and the k -packing independence number capturing upper bound and lower bound for the regret respectively. We show that the two notions are inherently connected via aligning them with the linear program of the weakly dominating set and its dual --- the fractional vertex packing set respectively. Based on this connection, we utilize the strong duality theorem to prove a general regret upper bound $O(\log |V| \cdot \delta \cdot \log |V|)^{\frac{1}{3}}$ and a lower bound $\Omega(\log |V| \cdot \alpha)^{\frac{1}{3}}$ where α is the integrality gap of the dual linear program. Therefore, our bounds are tight up to a $(\log |V|)^{\frac{1}{3}}$ factor on graphs with bounded integrality gap for the vertex packing problem including trees and graphs with bounded degree. Moreover, we show that for several special families of graphs, we can get rid of the $(\log |V|)^{\frac{1}{3}}$ factor and establish optimal regret.

Information-theoretic generalization bounds for black-box learning algorithms

- Hrayr Harutyunyan, Maxim Raginsky, Greg Ver Steeg, Aram Galstyan
- abstract@[open-review](#): We derive information-theoretic generalization bounds for supervised learning algorithms based on the information contained in predictions rather than in the output of the training algorithm. These bounds improve over the existing information-theoretic bounds, are applicable to a wider range of algorithms, and solve two key challenges: (a) they give meaningful results for deterministic algorithms and (b) they are significantly easier to estimate. We show experimentally that the proposed bounds closely follow the generalization gap in practical scenarios for deep learning.

Trash or Treasure? An Interactive Dual-Stream Strategy for Single Image Reflection Separation

- Qiming Hu, Xiaojie Guo
- abstract@[open-review](#): Single image reflection separation (SIRS), as a representative blind source separation task, aims to recover two layers, transmission and reflection, from one mixed observation, which is challenging due to the highly ill-posed nature. Existing deep learning based solutions typically restore the target layers individually, or with some concerns at the end of the output, barely taking into account the interaction across the two streams/branches. In order to utilize information more efficiently, this work presents a general yet simple interactive strategy, namely

\$textit{your trash is my treasure}\$(YTMT), for constructing dual-stream decomposition networks. To be specific, we explicitly enforce the two streams to communicate with each other block-wisely. Inspired by the additive property between the two components, the interactive path can be easily built via transferring, instead of discarding, deactivated information by the ReLU rectifier from one stream to the other. Both ablation studies and experimental results on widely-used SIRS datasets are conducted to demonstrate the efficacy of YTMT, and reveal its superiority over other state-of-the-art alternatives. The implementation is quite simple and our code is publicly available at <https://github.com/mingcv/YTMT-Strategy>.

[Rot-Pro: Modeling Transitivity by Projection in Knowledge Graph Embedding](#)

- Tengwei Song, Jie Luo, Lei Huang
- abstract@[open-review](#): Knowledge graph embedding models learn the representations of entities and relations in the knowledge graphs for predicting missing links (relations) between entities. Their effectiveness are deeply affected by the ability of modeling and inferring different relation patterns such as symmetry, asymmetry, inversion, composition and transitivity. Although existing models are already able to model many of these relations patterns, transitivity, a very common relation pattern, is still not been fully supported. In this paper, we first theoretically show that the transitive relations can be modeled with projections. We then propose the Rot-Pro model which combines the projection and relational rotation together. We prove that Rot-Pro can infer all the above relation patterns. Experimental results show that the proposed Rot-Pro model effectively learns the transitivity pattern and achieves the state-of-the-art results on the link prediction task in the datasets containing transitive relations.

[Planning from Pixels in Environments with Combinatorially Hard Search Spaces](#)

- Marco Bagatella, Miroslav Olšák, Michal Rolínek, Georg Martius
- abstract@[open-review](#): The ability to form complex plans based on raw visual input is a litmus test for current capabilities of artificial intelligence, as it requires a seamless combination of visual processing and abstract algorithmic execution, two traditionally separate areas of computer science. A recent surge of interest in this field brought advances that yield good performance in tasks ranging from arcade games to continuous control; these methods however do not come without significant issues, such as limited generalization capabilities and difficulties when dealing with combinatorially hard planning instances. Our contribution is two-fold: (i) we present a method that learns to represent its environment as a latent graph and leverages state reidentification to reduce the complexity of finding a good policy from exponential to linear (ii) we introduce a set of lightweight environments with an underlying discrete combinatorial structure in which planning is challenging even for humans. Moreover, we show that our methods achieves strong empirical generalization to variations in the environment, even across highly disadvantaged regimes, such as “one-shot” planning, or in an offline RL paradigm which only provides low-quality trajectories.

[PLUGIn: A simple algorithm for inverting generative models with recovery guarantees](#)

- Babhru Joshi, Xiaowei Li, Yaniv Plan, Ozgur Yilmaz
- abstract@[open-review](#): We consider the problem of recovering an unknown latent code vector under a known generative model. For a d -layer deep generative network $\mathcal{G}: \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_d}$ with ReLU activation functions, let the observation be $\mathcal{G}(x) + \epsilon$ where ϵ is noise. We introduce a simple novel algorithm, Partially Linearized Update for Generative Inversion (PLUGIn), to estimate x (and thus $\mathcal{G}(x)$). We prove that, when weights are Gaussian and layer widths $n_i \leq 5^i n_0$ (up to log factors), the algorithm converges geometrically to a neighbourhood of x with high probability. Note the inequality on layer widths allows $n_i > n_{i+1}$ when $i \geq 1$. To our knowledge, this is the first such result for networks with some contractive layers. After a sufficient number of iterations, the estimation errors for both x and $\mathcal{G}(x)$ are at most in the order of $\sqrt{4^n n_0 / n_d} \|\epsilon\|$. Thus, the algorithm can denoise when the expansion ratio n_d/n_0 is large. Numerical experiments on synthetic data and real data are provided to validate our theoretical results and to illustrate that the algorithm can effectively remove artifacts in an image.

[Modular Gaussian Processes for Transfer Learning](#)

- Pablo Moreno-Muñoz, Antonio Artes, Mauricio Álvarez
- abstract@[open-review](#): We present a framework for transfer learning based on modular variational Gaussian processes (GP). We develop a module-based method that having a dictionary of well fitted GPs, each model being characterised by its hyperparameters, pseudo-inputs and their corresponding posterior densities, one could build ensemble GP models without revisiting any data. Our method avoids undesired data centralisation, reduces rising computational costs and allows the transfer of learned uncertainty metrics after training. We exploit the augmentation of high-dimensional integral operators based on the Kullback-Leibler divergence between stochastic processes to introduce an efficient lower bound under all the sparse variational GPs, with different complexity and even likelihood distribution. The method is also valid for multi-output GPs, learning correlations a posteriori between independent modules. Extensive results illustrate the usability of our framework in large-scale and multi-task experiments, also compared with the exact inference methods in the literature.

[Neural Human Performer: Learning Generalizable Radiance Fields for Human Performance Rendering](#)

- Youngjoong Kwon, Dahun Kim, Duygu Ceylan, Henry Fuchs
- abstract@[open-review](#): In this paper, we aim at synthesizing a free-viewpoint video of an arbitrary human performance using sparse multi-view cameras. Recently, several works have addressed this problem by learning person-specific neural radiance fields (NeRF) to capture the appearance of a particular human. In parallel, some work proposed to use pixel-aligned features to generalize radiance fields to arbitrary new scenes and objects. Adopting such generalization approaches to humans, however, is highly challenging due to the heavy occlusions and dynamic articulations of body parts. To tackle this, we propose Neural Human Performer, a novel approach that learns generalizable neural radiance fields based on a parametric human body model for robust performance capture. Specifically, we first introduce a temporal transformer that aggregates tracked visual features based on the skeletal body motion over time. Moreover, a multi-view transformer is proposed to perform cross-attention between the temporally-fused features and the pixel-aligned features at each time step to integrate observations on the fly from multiple views. Experiments on the ZJU-MoCap and AIST datasets show that our method significantly outperforms recent generalizable NeRF methods on unseen identities and poses.

[Locally differentially private estimation of functionals of discrete distributions](#)

- Cristina Butucea, Yann ISSARTEL
- abstract@[open-review](#): We study the problem of estimating non-linear functionals of discrete distributions in the context of local differential privacy. The initial data $x_1, \dots, x_n \in [K]$ are supposed i.i.d. and distributed according to an unknown discrete distribution $p = (p_1, \dots, p_K)$. Only α -locally differentially private (LDP) samples z_1, \dots, z_n are publicly available, where the term ‘local’ means that each z_i is produced using one individual attribute x_i . We exhibit privacy mechanisms (PM) that are interactive (i.e. they are allowed to use already published confidential data) or non-interactive. We describe the behavior of the quadratic risk for estimating the power sum functional $F_\gamma = \sum_{k=1}^K p_k^\gamma$, $\gamma > 0$ as a function of K , n and α . In the non-interactive case, we study two plug-in type estimators of F_γ , for all $\gamma > 0$, that are similar to the MLE analyzed by Jiao et al. (2017) in the multinomial model. However, due to the privacy constraint the rates we attain are slower and similar to those obtained in the Gaussian model by Collier et al. (2020). In the sequentially interactive case, we introduce for all $\gamma > 1$ a two-step procedure which attains the parametric rate $(n \alpha^2)^{-1/2}$ when $\gamma \geq 2$. We give lower bounds results over all α -LDP mechanisms and over all estimators using the private samples.

[Asymptotics of representation learning in finite Bayesian neural networks](#)

- Jacob Zavatone-Veth, Abdulkadir Canatar, Ben Ruben, Cengiz Pehlevan
- abstract@[open-review](#): Recent works have suggested that finite Bayesian neural networks may sometimes outperform their infinite cousins because finite networks can flexibly adapt their internal representations. However, our theoretical understanding of how the learned hidden layer representations of finite networks differ from the fixed representations of infinite networks remains incomplete. Perturbative finite-width corrections to the network prior and posterior have been studied, but the asymptotics of learned features have not been fully characterized. Here, we argue that the leading finite-width corrections to the average feature kernels for any Bayesian network with linear readout and Gaussian likelihood have a largely universal form. We illustrate this explicitly for three tractable network architectures: deep linear fully-connected and convolutional networks, and networks with a single nonlinear hidden layer. Our results begin to elucidate how task-relevant learning signals shape the hidden layer representations of wide Bayesian neural networks.

[Adaptive Ensemble Q-learning: Minimizing Estimation Bias via Error Feedback](#)

- Hang Wang, Sen Lin, Junshan Zhang
- abstract@[open-review](#): The ensemble method is a promising way to mitigate the overestimation issue in Q-learning, where multiple function approximators are used to estimate the action values. It is known that the estimation bias hinges heavily on the ensemble size (i.e., the number of Q-function approximators used in the target), and that determining the 'right' ensemble size is highly nontrivial, because of the time-varying nature of the function approximation errors during the learning process. To tackle this challenge, we first derive an upper bound and a lower bound on the estimation bias, based on which the ensemble size is adapted to drive the bias to be nearly zero, thereby coping with the impact of the time-varying approximation errors accordingly. Motivated by the theoretic findings, we advocate that the ensemble method can be combined with Model Identification Adaptive Control (MIAC) for effective ensemble size adaptation. Specifically, we devise Adaptive Ensemble Q-learning (AdaEQ), a generalized ensemble method with two key steps: (a) approximation error characterization which serves as the feedback for flexibly controlling the ensemble size, and (b) ensemble size adaptation tailored towards minimizing the estimation bias. Extensive experiments are carried out to show that AdaEQ can improve the learning performance than the existing methods for the MuJoCo benchmark.

[Domain Adaptation with Invariant Representation Learning: What Transformations to Learn?](#)

- Petar Stojanov, Zijian Li, Mingming Gong, Ruichu Cai, Jaime Carbonell, Kun Zhang
- abstract@[open-review](#): Unsupervised domain adaptation, as a prevalent transfer learning setting, spans many real-world applications. With the increasing representational power and applicability of neural networks, state-of-the-art domain adaptation methods make use of deep architectures to map the input features $\$X\$$ to a latent representation $\$Z\$$ that has the same marginal distribution across domains. This has been shown to be insufficient for generating optimal representation for classification, and to find conditionally invariant representations, usually strong assumptions are needed. We provide reasoning why when the supports of the source and target data from overlap, any map of $\$X\$$ that is fixed across domains may not be suitable for domain adaptation via invariant features. Furthermore, we develop an efficient technique in which the optimal map from $\$X\$$ to $\$Z\$$ also takes domain-specific information as input, in addition to the features $\$X\$$. By using the property of minimal changes of causal mechanisms across domains, our model also takes into account the domain-specific information to ensure that the latent representation $\$Z\$$ does not discard valuable information about $\$Y\$$. We demonstrate the efficacy of our method via synthetic and real-world data experiments. The code is available at: \texttt{https://github.com/DMIRLAB-Group/DSAN}.

[CSDI: Conditional Score-based Diffusion Models for Probabilistic Time Series Imputation](#)

- Yusuke Tashiro, Jiaming Song, Yang Song, Stefano Ermon
- abstract@[open-review](#): The imputation of missing values in time series has many applications in healthcare and finance. While autoregressive models are natural candidates for time series imputation, score-based diffusion models have recently outperformed existing counterparts including autoregressive models in many tasks such as image generation and audio synthesis, and would be promising for time series imputation. In this paper, we propose Conditional Score-based Diffusion model (CSDI), a novel time series imputation method that utilizes score-based diffusion models conditioned on observed data. Unlike existing score-based approaches, the conditional diffusion model is explicitly trained for imputation and can exploit correlations between observed values. On healthcare and environmental data, CSDI improves by 40-65% over existing probabilistic imputation methods on popular performance metrics. In addition, deterministic imputation by CSDI reduces the error by 5-20% compared to the state-of-the-art deterministic imputation methods. Furthermore, CSDI can also be applied to time series interpolation and probabilistic forecasting, and is competitive with existing baselines. The code is available at <https://github.com/ermongroup/CSDI>.

[Causal Bandits with Unknown Graph Structure](#)

- Yangyi Lu, Amirhossein Meisami, Ambuj Tewari
- abstract@[open-review](#): In causal bandit problems the action set consists of interventions on variables of a causal graph. Several researchers have recently studied such bandit problems and pointed out their practical applications. However, all existing works rely on a restrictive and impractical assumption that the learner is given full knowledge of the causal graph structure upfront. In this paper, we develop novel causal bandit algorithms without knowing the causal graph. Our algorithms work well for causal trees, causal forests and a general class of causal graphs. The regret guarantees of our algorithms greatly improve upon those of standard multi-armed bandit (MAB) algorithms under mild conditions. Lastly, we prove our mild conditions are necessary: without them one cannot do better than standard MAB algorithms.

[Piper: Multidimensional Planner for DNN Parallelization](#)

- Jakub M. Tarnawski, Deepak Narayanan, Amar Phanishayee
- abstract@[open-review](#): The rapid increase in sizes of state-of-the-art DNN models, and consequently the increase in the compute and memory requirements of model training, has led to the development of many execution schemes such as data parallelism, pipeline model parallelism, tensor (intra-layer) model parallelism, and various memory-saving optimizations. However, no prior work has tackled the highly complex problem of optimally partitioning the DNN computation graph across many accelerators while combining all these parallelism modes and optimizations. In this work, we introduce Piper, an efficient optimization algorithm for this problem that is based on a two-level dynamic programming approach. Our two-level approach is driven by the insight that being given tensor-parallelization techniques for individual layers (e.g., Megatron-LM's splits for transformer layers) significantly reduces the search space and makes the global problem tractable, compared to considering tensor-parallel configurations for the entire DNN operator graph.

[Causal Effect Inference for Structured Treatments](#)

- Jean Kaddour, Yuchen Zhu, Qi Liu, Matt J. Kusner, Ricardo Silva
- abstract@[open-review](#): We address the estimation of conditional average treatment effects (CATEs) for structured treatments (e.g., graphs, images, texts). Given a weak condition on the effect, we propose the generalized Robinson decomposition, which (i) isolates the causal estimand (reducing regularization bias), (ii) allows one to plug in arbitrary models for learning, and (iii) possesses a quasi-oracle convergence guarantee under mild assumptions. In experiments with small-world and molecular graphs we demonstrate that our approach outperforms prior work in CATE estimation.

[Efficient hierarchical Bayesian inference for spatio-temporal regression models in neuroimaging](#)

- Ali Hashemi, Yijing Gao, Chang Cai, Sanjay Ghosh, Klaus-Robert Müller, Srikanth Nagarajan, Stefan Haufe
- abstract@[open-review](#): Several problems in neuroimaging and beyond require inference on the parameters of multi-task sparse hierarchical regression models. Examples include M/EEG inverse problems, neural encoding models for task-based fMRI analyses, and climate science. In these domains, both the model parameters to be inferred and the measurement noise may exhibit a complex spatio-temporal structure. Existing work either neglects the temporal structure or leads to computationally demanding inference schemes. Overcoming these limitations, we devise a novel flexible hierarchical Bayesian framework within which the spatio-temporal dynamics of model parameters and noise are modeled to have Kronecker product covariance structure. Inference in our framework is based on majorization-minimization optimization and has guaranteed convergence properties. Our highly efficient algorithms exploit the intrinsic Riemannian geometry of temporal autocovariance matrices. For stationary dynamics described by Toeplitz matrices, the theory of circulant embeddings is employed. We prove convex bounding properties and derive update rules of the resulting algorithms. On both synthetic and real neural data from M/EEG, we demonstrate that our methods lead to improved performance.

[Topological Attention for Time Series Forecasting](#)

- Sebastian Zeng, Florian Graf, Christoph Hofer, Roland Kwitt
- abstract@[open-review](#): The problem of (point) forecasting univariate time series is considered. Most approaches, ranging from traditional statistical methods to recent learning-based techniques with neural networks, directly operate on raw time series observations. As an extension, we study whether local topological properties, as captured via persistent homology, can serve as a reliable signal that provides complementary information for learning to forecast. To this end, we propose topological attention, which allows attending to local topological features within a time horizon of historical data. Our approach easily integrates into existing end-to-end trainable forecasting models, such as N-BEATS, and, in combination with the latter exhibits state-of-the-art performance on the large-scale M4 benchmark dataset of 100,000 diverse time series from different domains. Ablation experiments, as well as a comparison to recent techniques in a setting where only a single time series is available for training, corroborate the beneficial nature of including local topological information through an attention mechanism.

[Local Signal Adaptivity: Provable Feature Learning in Neural Networks Beyond Kernels](#)

- Stefani Karp, Ezra Winston, Yuanzhi Li, Aarti Singh
- abstract@[open-review](#): Neural networks have been shown to outperform kernel methods in practice (including neural tangent kernels). Most theoretical explanations of this performance gap focus on learning a complex hypothesis class; in some cases, it is unclear whether this hypothesis class captures realistic data. In this work, we propose a related, but alternative, explanation for this performance gap in the image classification setting, based on finding a sparse signal in the presence of noise. Specifically, we prove that, for a simple data distribution with sparse signal amidst high-variance noise, a simple convolutional neural network trained using stochastic gradient descent learns to threshold out the noise and find the signal. On the other hand, the corresponding neural tangent kernel, with a fixed set of predetermined features, is unable to adapt to the signal in this manner. We supplement our theoretical results by demonstrating this phenomenon empirically: in CIFAR-10 and MNIST images with various backgrounds, as the background noise increases in intensity, a CNN's performance stays relatively robust, whereas its corresponding neural tangent kernel sees a notable drop in performance. We therefore propose the "local signal adaptivity" (LSA) phenomenon as one explanation for the superiority of neural networks over kernel methods.

[IA-RED\\$^2\\$: Interpretability-Aware Redundancy Reduction for Vision Transformers](#)

- Bowen Pan, Rameswar Panda, Yifan Jiang, Zhangyang Wang, Rogerio Feris, Aude Oliva
- abstract@[open-review](#): The self-attention-based model, transformer, is recently becoming the leading backbone in the field of computer vision. In spite of the impressive success made by transformers in a variety of vision tasks, it still suffers from heavy computation and intensive memory costs. To address this limitation, this paper presents an Interpretability-Aware REDundancy REDuction framework (IA-RED\$^2\$). We start by observing a large amount of redundant computation, mainly spent on uncorrelated input patches, and then introduce an interpretable module to dynamically and gracefully drop these redundant patches. This novel framework is then extended to a hierarchical structure, where uncorrelated tokens at different stages are gradually removed, resulting in a considerable shrinkage of computational cost. We include extensive experiments on both image and video tasks, where our method could deliver up to 1.4x speed-up for state-of-the-art models like DeiT and TimeSformer, by only sacrificing less than 0.7% accuracy. More importantly, contrary to other acceleration approaches, our method is inherently interpretable with substantial visual evidence, making vision transformer closer to a more human-understandable architecture while being lighter. We demonstrate that the interpretability that naturally emerged in our framework can outperform the raw attention learned by the original visual transformer, as well as those generated by off-the-shelf interpretation methods, with both qualitative and quantitative results. Project Page: <http://people.csail.mit.edu/bpan/ia-red/>.

[Symbolic Regression via Deep Reinforcement Learning Enhanced Genetic Programming Seeding](#)

- Terrell Mundhenk, Mikel Landajuela, Ruben Glatt, Claudio P Santiago, Daniel faissol, Brenden K Petersen
- abstract@[open-review](#): Symbolic regression is the process of identifying mathematical expressions that fit observed output from a black-box process. It is a discrete optimization problem generally believed to be NP-hard. Prior approaches to solving the problem include neural-guided search (e.g. using reinforcement learning) and genetic programming. In this work, we introduce a hybrid neural-guided/genetic programming approach to symbolic regression and other combinatorial optimization problems. We propose a neural-guided component used to seed the starting population of a random restart genetic programming component, gradually learning better starting populations. On a number of common benchmark tasks to recover underlying expressions from a dataset, our method recovers 65% more expressions than a recently published top-performing model using the same experimental setup. We demonstrate that running many genetic programming generations without interdependence on the neural-guided component performs better for symbolic regression than alternative formulations where the two are more strongly coupled. Finally, we introduce a new set of 22 symbolic regression benchmark problems with increased difficulty over existing benchmarks. Source code is provided at www.github.com/brendenpetersen/deep-symbolic-optimization.

[Choose a Transformer: Fourier or Galerkin](#)

- Shuhao Cao
- abstract@[open-review](#): In this paper, we apply the self-attention from the state-of-the-art Transformer in Attention Is All You Need for the first time to a data-driven operator learning problem related to partial differential equations. An effort is put together to explain the heuristics of, and to improve the efficacy of the attention mechanism. By employing the operator approximation theory in Hilbert spaces, it is demonstrated for the first time that the softmax normalization in the scaled dot-product attention is sufficient but not necessary. Without softmax, the approximation capacity of a linearized Transformer variant can be proved to be comparable to a Petrov-Galerkin projection layer-wise, and the estimate is independent with respect to the sequence length. A new layer normalization scheme mimicking the Petrov-Galerkin projection is proposed to allow a scaling to propagate through attention layers, which helps the model achieve remarkable accuracy in operator learning tasks with unnormalized data. Finally, we present three operator learning experiments, including the viscous Burgers' equation, an interface Darcy flow, and an inverse interface coefficient identification problem. The newly proposed simple attention-based operator learner, Galerkin Transformer, shows significant improvements in both training cost and evaluation accuracy over its softmax-normalized counterparts.

[A Causal Lens for Controllable Text Generation](#)

- Zhiting Hu, Li Erran Li
- abstract@[open-review](#): Controllable text generation concerns two fundamental tasks of wide applications, namely generating text of given attributes (i.e., attribute-conditional generation), and minimally editing existing text to possess desired attributes (i.e., text attribute transfer). Extensive prior work has largely studied the two problems separately, and developed different conditional models which, however, are prone to producing biased text (e.g., various gender stereotypes). This paper proposes to formulate controllable text generation from a principled causal perspective which models the two tasks with a unified framework. A direct advantage of the causal formulation is the use of rich causality tools to mitigate generation biases and improve control. We treat the two tasks as interventional and counterfactual causal inference based on a structural causal model, respectively. We then apply the framework to the challenging practical setting where confounding factors (that induce spurious correlations) are observable only on a small fraction of data. Experiments show significant superiority of the causal approach over previous conditional models for improved control accuracy and reduced bias.

[Differentially Private Multi-Armed Bandits in the Shuffle Model](#)

- Jay Tenenbaum, Haim Kaplan, Yishay Mansour, Uri Stemmer
- abstract@[open-review](#): We give an (ε, δ) -differentially private algorithm for the Multi-Armed Bandit (MAB) problem in the shuffle model with a distribution-dependent regret of $O(\left(\sum_{a: \Delta_a > 0} \frac{\log T}{\Delta_a}\right) + \frac{k \sqrt{\log \frac{1}{\delta}}}{\Delta_a} \log T)$, and a distribution-independent regret of $O(\sqrt{kT \log T} + \frac{k \sqrt{\log \frac{1}{\delta}}}{\Delta_a} \log T)$, where T is the number of rounds, Δ_a is the suboptimality gap of the action a , and k is the total number of actions. Our upper bound almost matches the regret of the best known algorithms for the centralized model, and significantly outperforms the best known algorithm in the local model.

[Dual Adaptivity: A Universal Algorithm for Minimizing the Adaptive Regret of Convex Functions](#)

- Lijun Zhang, Guanghui Wang, Wei-Wei Tu, Wei Jiang, Zhi-Hua Zhou
- abstract@[open-review](#): To deal with changing environments, a new performance measure—adaptive regret, defined as the maximum static regret over any interval, was proposed in online learning. Under the setting of online convex optimization, several algorithms have been successfully developed to minimize the adaptive regret. However, existing algorithms lack universality in the sense that they can only handle one type of convex functions and need apriori knowledge of parameters. By contrast, there exist universal algorithms, such as MetaGrad, that attain optimal static regret for multiple types of convex functions simultaneously. Along this line of research, this paper presents the first universal algorithm for minimizing the adaptive regret of convex functions. Specifically, we borrow the idea of maintaining multiple learning rates in MetaGrad to handle the uncertainty of functions, and utilize the technique of sleeping experts to capture changing environments. In this way, our algorithm automatically adapts to the property of functions (convex, exponentially concave, or strongly convex), as well as the nature of environments (stationary or changing). As a by product, it also allows the type of functions to switch between rounds.

[Learning Hard Optimization Problems: A Data Generation Perspective](#)

- James Kotary, Ferdinando Fioretto, Pascal Van Hentenryck
- abstract@[open-review](#): Optimization problems are ubiquitous in our societies and are present in almost every segment of the economy. Most of these optimization problems are NP-hard and computationally demanding, often requiring approximate solutions for large-scale instances. Machine learning frameworks that learn to approximate solutions to such hard optimization problems are a potentially promising avenue to address these difficulties, particularly when many closely related problem instances must be solved repeatedly. Supervised learning frameworks can train a model using the outputs of pre-solved instances. However, when the outputs are themselves approximations, when the optimization problem has symmetric solutions, and/or when the solver uses randomization, solutions to closely related instances may exhibit large differences and the learning task can become inherently more difficult. This paper demonstrates this critical challenge, connects the volatility of the training data to the ability of a model to approximate it, and proposes a method for producing (exact or approximate) solutions to optimization problems that are more amenable to supervised learning tasks. The effectiveness of the method is tested on hard non-linear nonconvex and discrete combinatorial problems.

[Canonical Capsules: Self-Supervised Capsules in Canonical Pose](#)

- Weiwei Sun, Andrea Tagliasacchi, Boyang Deng, Sara Sabour, Soroosh Yazdani, Geoffrey E. Hinton, Kwang Moo Yi
- abstract@[open-review](#): We propose a self-supervised capsule architecture for 3D point clouds. We compute capsule decompositions of objects through permutation-equivariant attention, and self-supervise the process by training with pairs of randomly rotated objects. Our key idea is to aggregate the attention masks into semantic keypoints, and use these to supervise a decomposition that satisfies the capsule invariance/equivariance properties. This not only enables the training of a semantically consistent decomposition, but also allows us to learn a canonicalization operation that enables object-centric reasoning. To train our neural network we require neither classification labels nor manually-aligned training datasets. Yet, by learning an object-centric representation in a self-supervised manner, our method outperforms the state-of-the-art on 3D point cloud reconstruction, canonicalization, and unsupervised classification.

[Characterizing Generalization under Out-Of-Distribution Shifts in Deep Metric Learning](#)

- Timo Milbich, Karsten Roth, Samarth Sinha, Ludwig Schmidt, Marzyeh Ghassemi, Bjorn Ommer
- abstract@[open-review](#): Deep Metric Learning (DML) aims to find representations suitable for zero-shot transfer to a priori unknown test distributions. However, common evaluation protocols only test a single, fixed data split in which train and test classes are assigned randomly. More realistic evaluations should consider a broad spectrum of distribution shifts with potentially varying degree and difficulty. In this work, we systematically construct train-test splits of increasing difficulty and present the ooDML benchmark to characterize generalization under out-of-distribution shifts in DML. ooDML is designed to probe the generalization performance on much more challenging, diverse train-to-test distribution shifts. Based on our new benchmark, we conduct a thorough empirical analysis of state-of-the-art DML methods. We find that while generalization tends to consistently degrade with difficulty, some methods are better at retaining performance as the distribution shift increases. Finally, we propose few-shot DML as an efficient way to consistently improve generalization in response to unknown test shifts presented in ooDML.

[Dynamics-regulated kinematic policy for egocentric pose estimation](#)

- Zhengyi Luo, Ryo Hachiuma, Ye Yuan, Kris Kitani
- abstract@[open-review](#): We propose a method for object-aware 3D egocentric pose estimation that tightly integrates kinematics modeling, dynamics modeling, and scene object information. Unlike prior kinematics or dynamics-based approaches where the two components are used disjointly, we synergize the two approaches via dynamics-regulated training. At each timestep, a kinematic model is used to provide a target pose using video evidence and simulation state. Then, a prelearned dynamics model attempts to mimic the kinematic pose in a physics simulator. By comparing the pose instructed by the kinematic model against the pose generated by the dynamics model, we can use their misalignment to further improve the kinematic model. By factoring in the 6DoF pose of objects (e.g., chairs, boxes) in the scene, we demonstrate for the first time, the ability to estimate physically-plausible 3D human-object interactions using a single wearable camera. We evaluate our egocentric pose estimation method in both controlled laboratory settings and real-world scenarios.

[Never Go Full Batch \(in Stochastic Convex Optimization\)](#)

- Idan Amir, Yair Carmon, Tomer Koren, Roi Livni
- abstract@[open-review](#): We study the generalization performance of \$text{\emph{full-batch}}\$ optimization algorithms for stochastic convex optimization: these are first-order methods that only access the exact gradient of the empirical risk (rather than gradients with respect to individual data points), that include a wide range of algorithms such as gradient descent, mirror descent, and their regularized and/or accelerated variants. We provide a new separation result showing that, while algorithms such as stochastic gradient descent can generalize and optimize the population risk to within \$\epsilon\$ after \$O(1/\epsilon^2)\$ iterations, full-batch methods either need at least \$\Omega(1/\epsilon^4)\$ iterations or exhibit a dimension-dependent sample complexity.

[Collaborative Learning in the Jungle \(Decentralized, Byzantine, Heterogeneous, Asynchronous and Nonconvex Learning\)](#)

- El Mahdi El-Mhamdi, Sadegh Farhadkhani, Rachid Guerraoui, Arsany Guirguis, Lê-Nguyễn Hoang, Sébastien Rouault
- abstract@[open-review](#): We study \emph{Byzantine collaborative learning}, where \$n\$ nodes seek to collectively learn from each others' local data. The data distribution may vary from one node to another. No node is trusted, and \$f < n\$ nodes can behave arbitrarily. We prove that collaborative learning is equivalent to a new form of agreement, which we call \emph{averaging agreement}. In this problem, nodes start each with an initial vector and seek to approximately agree on a common vector, which is close to the average of honest nodes' initial vectors. We present two asynchronous solutions to averaging agreement, each we prove optimal according to some dimension. The first, based on the minimum-diameter averaging, requires \$n \geq 6f+1\$, but achieves asymptotically the best-possible averaging constant up to a multiplicative constant. The second, based on reliable broadcast and coordinate-wise trimmed mean, achieves optimal Byzantine resilience, i.e., \$n \geq 3f+1\$. Each of these algorithms induces an optimal Byzantine collaborative learning protocol. In particular, our equivalence yields new impossibility theorems on what any collaborative learning algorithm can achieve in adversarial and heterogeneous environments.

[Not All Low-Pass Filters are Robust in Graph Convolutional Networks](#)

- Heng Chang, Yu Rong, Tingyang Xu, Yatao Bian, Shiji Zhou, Xin Wang, Junzhou Huang, Wenwu Zhu
- abstract@[open-review](#): Graph Convolutional Networks (GCNs) are promising deep learning approaches in learning representations for graph-structured data. Despite the proliferation of such methods, it is well known that they are vulnerable to carefully crafted adversarial attacks on the graph structure. In this paper, we first conduct an adversarial vulnerability analysis based on matrix perturbation theory. We prove that the low- frequency components of the symmetric normalized Laplacian, which is usually used as the convolutional filter in GCNs, could be more robust against structural perturbations when their eigenvalues fall into a certain robust interval. Our results indicate that not all low-frequency components are robust to adversarial attacks and provide a deeper understanding of the relationship between graph spectrum and robustness of GCNs. Motivated by the theory, we present GCN-LFR, a general robust co-training paradigm for GCN-based models, that encourages transferring the robustness of low-frequency components with an auxiliary neural network. To this end, GCN-LFR could enhance the robustness of various kinds of GCN-based models against poisoning structural attacks in a plug-and-play manner. Extensive experiments across five benchmark datasets and five GCN-based models also confirm that GCN-LFR is resistant to the adversarial attacks without compromising on performance in the benign situation.

[Counterfactual Maximum Likelihood Estimation for Training Deep Networks](#)

- Xinyi Wang, Wenhui Chen, Michael Saxon, William Yang Wang
- abstract@[open-review](#): Although deep learning models have driven state-of-the-art performance on a wide array of tasks, they are prone to spurious correlations that should not be learned as predictive clues. To mitigate this problem, we propose a causality-based training framework to reduce the spurious correlations caused by observed confounders. We give theoretical analysis on the underlying general Structural Causal Model (SCM) and propose to perform Maximum Likelihood Estimation (MLE) on the interventional distribution instead of the observational distribution, namely Counterfactual Maximum Likelihood Estimation (CMLE). As the interventional distribution, in general, is hidden from the observational data, we then derive two different upper bounds of the expected negative log-likelihood and propose two general algorithms, Implicit CMLE and Explicit CMLE, for causal predictions of deep learning models using observational data. We conduct experiments on both simulated data and two real-world tasks: Natural Language Inference (NLI) and Image Captioning. The results show that CMLE methods outperform the regular MLE method in terms of out-of-domain generalization performance and reducing spurious correlations, while maintaining comparable performance on the regular evaluations.

[Robust Optimization for Multilingual Translation with Imbalanced Data](#)

- Xian Li, Hongyu Gong
- abstract@[open-review](#): Multilingual models are parameter-efficient and especially effective in improving low-resource languages by leveraging crosslingual transfer. Despite recent advance in massive multilingual translation with ever-growing model and data, how to effectively train multilingual models has not been well understood. In this paper, we show that a common situation in multilingual training, data imbalance among languages, poses optimization tension between high resource and low resource languages where the found multilingual solution is often sub-optimal for low resources. We show that common training method which upsamples low resources can not robustly optimize population loss with risks of either underfitting high resource languages or overfitting low resource ones. Drawing on recent findings on the geometry of loss landscape and its effect on generalization, we propose a principled optimization algorithm, Curvature Aware Task Scaling (CATS), which adaptively rescales gradients from different tasks with a meta objective of guiding multilingual training to low-curvature neighborhoods with uniformly low loss for all languages. We ran experiments on common benchmarks (TED, WMT and OPUS-100) with varying degrees of data imbalance. CATS effectively improved multilingual optimization and as a result demonstrated consistent gains on low resources (\$+0.8\$ to \$+2.2\$ BLEU) without hurting high resources. In addition, CATS is robust to overparameterization and large batch size training, making it a promising training method for massive multilingual models that truly improve low resource languages.

[A/B/n Testing with Control in the Presence of Subpopulations](#)

- Yoan Russac, Christina Katsimerou, Dennis Bohle, Olivier Cappé, Aurélien Garivier, Wouter M. Koolen
- abstract@[open-review](#): Motivated by A/B/n testing applications, we consider a finite set of distributions (called \emph{arms}), one of which is treated as a \emph{control}. We assume that the population is stratified into homogeneous subpopulations. At every time step, a subpopulation is sampled and an arm is chosen: the resulting observation is an independent draw from the arm conditioned on the subpopulation. The quality of each arm is assessed through a weighted combination of its subpopulation means. We propose a strategy for sequentially choosing one arm per time step so as to discover as fast as possible which arms, if any, have higher weighted expectation than the control. This strategy is shown to be asymptotically optimal in the following sense: if \$\tau_\delta\$ is the first time when the strategy ensures that it is able to output the correct answer with probability at least \$1-\delta\$, then \$\mathbb{E}[\tau_\delta]\$ grows linearly with \$\log(1/\delta)\$ at the exact optimal rate. This rate is identified in the paper in three different settings: (1) when the experimenter does not observe the subpopulation information, (2) when the subpopulation of each sample is observed but not chosen, and (3) when the experimenter can select the subpopulation from which each response is sampled. We illustrate the efficiency of the proposed strategy with numerical simulations on synthetic and real data collected from an A/B/n experiment.

[Using Random Effects to Account for High-Cardinality Categorical Features and Repeated Measures in Deep Neural Networks](#)

- Giora Simchoni, Saharon Rosset

- abstract@[open-review](#): High-cardinality categorical features are a major challenge for machine learning methods in general and for deep learning in particular. Existing solutions such as one-hot encoding and entity embeddings can be hard to scale when the cardinality is very high, require much space, are hard to interpret or may overfit the data. A special scenario of interest is that of repeated measures, where the categorical feature is the identity of the individual or object, and each object is measured several times, possibly under different conditions (values of the other features). We propose accounting for high-cardinality categorical features as random effects variables in a regression setting, and consequently adopt the corresponding negative log likelihood loss from the linear mixed models (LMM) statistical literature and integrate it in a deep learning framework. We test our model which we call LMMNN on simulated as well as real datasets with a single categorical feature with high cardinality, using various baseline neural networks architectures such as convolutional networks and LSTM, and various applications in e-commerce, healthcare and computer vision. Our results show that treating high-cardinality categorical features as random effects leads to a significant improvement in prediction performance compared to state of the art alternatives. Potential extensions such as accounting for multiple categorical features and classification settings are discussed. Our code and simulations are available at <https://github.com/gsimchoni/lmmnn>.

[Learning Debiased Representation via Disentangled Feature Augmentation](#)

- Jungsoo Lee, Eungyeup Kim, Juyoung Lee, Jihyeon Lee, Jaegul Choo
- abstract@[open-review](#): Image classification models tend to make decisions based on peripheral attributes of data items that have strong correlation with a target variable (i.e., dataset bias). These biased models suffer from the poor generalization capability when evaluated on unbiased datasets. Existing approaches for debiasing often identify and emphasize those samples with no such correlation (i.e., bias-conflicting) without defining the bias type in advance. However, such bias-conflicting samples are significantly scarce in biased datasets, limiting the debiasing capability of these approaches. This paper first presents an empirical analysis revealing that training with "diverse" bias-conflicting samples beyond a given training set is crucial for debiasing as well as the generalization capability. Based on this observation, we propose a novel feature-level data augmentation technique in order to synthesize diverse bias-conflicting samples. To this end, our method learns the disentangled representation of (1) the intrinsic attributes (i.e., those inherently defining a certain class) and (2) bias attributes (i.e., peripheral attributes causing the bias), from a large number of bias-aligned samples, the bias attributes of which have strong correlation with the target variable. Using the disentangled representation, we synthesize bias-conflicting samples that contain the diverse intrinsic attributes of bias-aligned samples by swapping their latent features. By utilizing these diversified bias-conflicting features during the training, our approach achieves superior classification accuracy and debiasing results against the existing baselines on both synthetic and real-world datasets.

[Scallop: From Probabilistic Deductive Databases to Scalable Differentiable Reasoning](#)

- Jian Huang, Ziyang Li, Binghong Chen, Karan Samel, Mayur Naik, Le Song, Xujie Si
- abstract@[open-review](#): Deep learning and symbolic reasoning are complementary techniques for an intelligent system. However, principled combinations of these techniques have limited scalability, rendering them ill-suited for real-world applications. We propose Scallop, a system that builds upon probabilistic deductive databases, to bridge this gap. The key insight underlying Scallop is a provenance framework that introduces a tunable parameter to specify the level of reasoning granularity. Scallop thereby i) generalizes exact probabilistic reasoning, ii) asymptotically reduces computational cost, and iii) provides relative accuracy guarantees. On a suite of tasks that involve mathematical and logical reasoning, Scallop scales significantly better without sacrificing accuracy compared to DeepProbLog, a principled neural logic programming approach. We also create and evaluate on a real-world Visual Question Answering (VQA) benchmark that requires multi-hop reasoning. Scallop outperforms two VQA-tailored models, a Neural Module Networks based and a transformer based model, by 12.42% and 21.66% respectively.

[Learning to Synthesize Programs as Interpretable and Generalizable Policies](#)

- Dweep Trivedi, Jesse Zhang, Shao-Hua Sun, Joseph J. Lim
- abstract@[open-review](#): Recently, deep reinforcement learning (DRL) methods have achieved impressive performance on tasks in a variety of domains. However, neural network policies produced with DRL methods are not human-interpretable and often have difficulty generalizing to novel scenarios. To address these issues, prior works explore learning programmatic policies that are more interpretable and structured for generalization. Yet, these works either employ limited policy representations (e.g. decision trees, state machines, or predefined program templates) or require stronger supervision (e.g. input/output state pairs or expert demonstrations). We present a framework that instead learns to synthesize a program, which details the procedure to solve a task in a flexible and expressive manner, solely from reward signals. To alleviate the difficulty of learning to compose programs to induce the desired agent behavior from scratch, we propose to first learn a program embedding space that continuously parameterizes diverse behaviors in an unsupervised manner and then search over the learned program embedding space to yield a program that maximizes the return for a given task. Experimental results demonstrate that the proposed framework not only learns to reliably synthesize task-solving programs but also outperforms DRL and program synthesis baselines while producing interpretable and more generalizable policies. We also justify the necessity of the proposed two-stage learning scheme as well as analyze various methods for learning the program embedding. Website at <https://clvrai.com/leaps>.

[The functional specialization of visual cortex emerges from training parallel pathways with self-supervised predictive learning](#)

- Shahab Bakhtiari, Patrick Mineault, Timothy Lillicrap, Christopher Pack, Blake Richards
- abstract@[open-review](#): The visual system of mammals is comprised of parallel, hierarchical specialized pathways. Different pathways are specialized in so far as they use representations that are more suitable for supporting specific downstream behaviours. In particular, the clearest example is the specialization of the ventral ("what") and dorsal ("where") pathways of the visual cortex. These two pathways support behaviours related to visual recognition and movement, respectively. To-date, deep neural networks have mostly been used as models of the ventral, recognition pathway. However, it is unknown whether both pathways can be modelled with a single deep ANN. Here, we ask whether a single model with a single loss function can capture the properties of both the ventral and the dorsal pathways. We explore this question using data from mice, who like other mammals, have specialized pathways that appear to support recognition and movement behaviours. We show that when we train a deep neural network architecture with two parallel pathways using a self-supervised predictive loss function, we can outperform other models in fitting mouse visual cortex. Moreover, we can model both the dorsal and ventral pathways. These results demonstrate that a self-supervised predictive learning approach applied to parallel pathway architectures can account for some of the functional specialization seen in mammalian visual systems.

[Adversarial Training Helps Transfer Learning via Better Representations](#)

- Zhun Deng, Linjun Zhang, Kailas Vodrahalli, Kenji Kawaguchi, James Y. Zou
- abstract@[open-review](#): Transfer learning aims to leverage models pre-trained on source data to efficiently adapt to target setting, where only limited data are available for model fine-tuning. Recent works empirically demonstrate that adversarial training in the source data can improve the ability of models to transfer to new domains. However, why this happens is not known. In this paper, we provide a theoretical model to rigorously analyze how adversarial training helps transfer learning. We show that adversarial training in the source data generates provably better representations, so fine-tuning on top of this representation leads to a more accurate predictor of the target data. We further demonstrate both theoretically and empirically that semi-supervised learning in the source data can also improve transfer learning by similarly improving the representation. Moreover, performing adversarial training on top of semi-supervised learning can further improve transferability, suggesting that the two approaches have complementary benefits on representations. We support our theories with experiments on popular data sets and deep learning architectures.

[Improving Coherence and Consistency in Neural Sequence Models with Dual-System, Neuro-Symbolic Reasoning](#)

- Maxwell Nye, Michael Tessler, Josh Tenenbaum, Brenden M. Lake
- abstract@[open-review](#): Human reasoning can be understood as an interplay between two systems: the intuitive and associative ("System 1") and the deliberative and logical ("System 2"). Neural sequence models---which have been increasingly successful at performing complex, structured tasks---exhibit the advantages and failure modes of System 1: they are fast and learn patterns from data, but are often inconsistent and incoherent. In this work, we seek a lightweight, training-free means of improving existing System 1-like sequence models by adding System 2-inspired logical reasoning. We explore several variations on this theme in which candidate generations from a neural sequence model are examined for logical consistency by a symbolic reasoning module, which can either accept or reject the generations. Our approach uses neural inference to mediate between the neural System 1 and the logical System 2. Results in robust story generation and grounded instruction-following show that this approach can increase the coherence and accuracy of neurally-based generations.

[Learning the optimal Tikhonov regularizer for inverse problems](#)

- Giovanni S. Alberti, Ernesto De Vito, Matti Lassas, Luca Ratti, Matteo Santacesaria
- abstract@[open-review](#): In this work, we consider the linear inverse problem $y = Ax + \varepsilon$, where $A : X \rightarrow Y$ is a known linear operator between the separable Hilbert spaces X and Y , x is a random variable in X and ε is a zero-mean random process in Y . This setting covers several inverse problems in imaging including denoising, deblurring, and X-ray tomography. Within the classical framework of regularization, we focus on the case where the regularization functional is not given a priori, but learned from data. Our first result is a characterization of the optimal generalized Tikhonov regularizer, with respect to the mean squared error. We find that it is completely independent of the forward operator A and depends only on the mean and covariance of x . Then, we consider the problem of learning the regularizer from a finite training set in two different frameworks: one supervised, based on samples of both x and y , and one unsupervised, based only on samples of x . In both cases, we prove generalization bounds, under some weak assumptions on the distribution of x and ε , including the case of sub-Gaussian variables. Our bounds hold in infinite-dimensional spaces, thereby showing that finer and finer discretizations do not make this learning problem harder. The results are validated through numerical simulations.

[NovelD: A Simple yet Effective Exploration Criterion](#)

- Tianjun Zhang, Huazhe Xu, Xiaolong Wang, Yi Wu, Kurt Keutzer, Joseph E. Gonzalez, Yuandong Tian
- abstract@[open-review](#): Efficient exploration under sparse rewards remains a key challenge in deep reinforcement learning. Previous exploration methods (e.g., RND) have achieved strong results in multiple hard tasks. However, if there are multiple novel areas to explore, these methods often focus quickly on one without sufficiently trying others (like a depth-wise first search manner). In some scenarios (e.g., four corridor environment in Sec 4.2), we observe they explore in one corridor for long and fail to cover all the states. On the other hand, in theoretical RL, with optimistic initialization and the inverse square root of visitation count as a bonus, it won't suffer from this and explores different novel regions alternatively (like a breadth-first search manner). In this paper, inspired by this, we propose a simple but effective criterion called NovelD by weighting every novel area approximately equally. Our algorithm is very simple but yet shows comparable performance or even outperforms multiple SOTA exploration methods in many hard exploration tasks. Specifically, NovelD solves all the static procedurally-generated tasks in Mini-Grid with just 120M environment steps, without any curriculum learning. In comparison, the previous SOTA only solves 50% of them. NovelD also achieves SOTA on multiple tasks in NetHack, a rogue-like game that contains more challenging procedurally-generated environments. In multiple Atari games (e.g., Montezuma's Revenge, Venture, Gravitar), NovelD outperforms RND. We analyze NovelD thoroughly in MiniGrid and found that empirically it helps the agent explore the environment more uniformly with a focus on exploring beyond the boundary.

[On Margin-Based Cluster Recovery with Oracle Queries](#)

- Marco Bressan, Nicolò Cesa-Bianchi, Silvio Lattanzi, Andrea Paudice
- abstract@[open-review](#): We study an active cluster recovery problem where, given a set of n points and an oracle answering queries like ``are these two points in the same cluster?'', the task is to recover exactly all clusters using as few queries as possible. We begin by introducing a simple but general notion of margin between clusters that captures, as special cases, the margins used in previous works, the classic SVM margin, and standard notions of stability for center-based clusterings. Under our margin assumptions we design algorithms that, in a variety of settings, recover all clusters exactly using only $O(\log n)$ queries. For R^m , we give an algorithm that recovers arbitrary convex clusters, in polynomial time, and with a number of queries that is lower than the best existing algorithm by $\Theta(m^m)$ factors. For general pseudometric spaces, where clusters might not be convex or might not have any notion of shape, we give an algorithm that achieves the $O(\log n)$ query bound, and is provably near-optimal as a function of the packing number of the space. Finally, for clusterings realized by binary concept classes, we give a combinatorial characterization of recoverability with $O(\log n)$ queries, and we show that, for many concept classes in R^m , this characterization is equivalent to our margin condition. Our results show a deep connection between cluster margins and active cluster recoverability.

[Multi-Scale Representation Learning on Proteins](#)

- Vignesh Ram Somnath, Charlotte Bunne, Andreas Krause
- abstract@[open-review](#): Proteins are fundamental biological entities mediating key roles in cellular function and disease. This paper introduces a multi-scale graph construction of a protein –HoloProt– connecting surface to structure and sequence. The surface captures coarser details of the protein, while sequence as primary component and structure –comprising secondary and tertiary components– capture finer details. Our graph encoder then learns a multi-scale representation by allowing each level to integrate the encoding from level(s) below with the graph at that level. We test the learned representation on different tasks, (i.) ligand binding affinity (regression), and (ii.) protein function prediction (classification). On the regression task, contrary to previous methods, our model performs consistently and reliably across different dataset splits, outperforming all baselines on most splits. On the classification task, it achieves a performance close to the top-performing model while using 10x fewer parameters. To improve the memory efficiency of our construction, we segment the multiplex protein surface manifold into molecular superpixels and substitute the surface with these superpixels at little to no performance loss.

[Sparse Quadratic Optimisation over the Stiefel Manifold with Application to Permutation Synchronisation](#)

- Florian Bernard, Daniel Cremers, Johan Thunberg
- abstract@[open-review](#): We address the non-convex optimisation problem of finding a sparse matrix on the Stiefel manifold (matrices with mutually orthogonal columns of unit length) that maximises (or minimises) a quadratic objective function. Optimisation problems on the Stiefel manifold occur for example in spectral relaxations of various combinatorial problems, such as graph matching, clustering, or permutation synchronisation. Although sparsity is a desirable property in such settings, it is mostly neglected in spectral formulations since existing solvers, e.g. based on eigenvalue decomposition, are unable to account for sparsity while at the same time maintaining global optimality guarantees. We fill this gap and propose a simple yet effective sparsity-promoting modification of the Orthogonal Iteration algorithm for finding the dominant eigenspace of a matrix. By doing so, we can guarantee that our method finds a Stiefel matrix that is globally optimal with respect to the quadratic objective function, while in addition being sparse. As a motivating application we consider the task of permutation synchronisation, which can be understood as a constrained clustering problem that has particular relevance for matching multiple images or 3D shapes in computer vision, computer graphics, and beyond. We demonstrate that the proposed approach outperforms previous methods in this domain.

[Second-Order Neural ODE Optimizer](#)

- Guan-Horng Liu, Tianrong Chen, Evangelos Theodorou
- abstract@[open-review](#): We propose a novel second-order optimization framework for training the emerging deep continuous-time models, specifically the Neural Ordinary Differential Equations (Neural ODEs). Since their training already involves expensive gradient computation by solving a backward ODE, deriving efficient second-order methods becomes highly nontrivial. Nevertheless, inspired by the recent Optimal Control (OC) interpretation of training deep networks, we show that a specific continuous-time OC methodology, called Differential Programming, can be adopted to derive backward ODEs for higher-order derivatives at the same $O(1)$ memory cost. We further explore a low-rank representation of the second-order derivatives and show that it leads to efficient preconditioned updates with the aid of Kronecker-based factorization. The resulting method – named SNOpt – converges much faster than first-order baselines in wall-clock time, and the improvement remains consistent across various applications, e.g. image classification, generative flow, and time-series prediction. Our framework also enables direct architecture optimization, such as the integration time of Neural ODEs, with second-order feedback policies, strengthening the OC perspective as a principled tool of analyzing optimization in deep learning. Our code is available at <https://github.com/ghliu/snopt>.

[Graph Neural Networks with Local Graph Parameters](#)

- Pablo Barceló, Floris Geerts, Juan Reutter, Maksimilian Ryschkov
- abstract@[open-review](#): Various recent proposals increase the distinguishing power of Graph Neural Networks (GNNs) by propagating features between k -tuples of vertices. The distinguishing power of these “higher-order” GNNs is known to be bounded by the k -dimensional Weisfeiler-Leman (WL) test, yet their $O(n^k)$ memory requirements limit their applicability. Other proposals infuse GNNs with local higher-order graph structural information from the start, hereby inheriting the desirable $O(n)$ memory requirement from GNNs at the cost of a one-time, possibly non-linear, preprocessing step. We propose local graph parameter enabled GNNs as a framework for studying the latter kind of approaches and precisely characterize their distinguishing power, in terms of a variant of the WL test, and in terms of the graph structural properties that they can take into account. Local graph parameters can be added to any GNN architecture, and are cheap to compute. In terms of expressive power, our proposal lies in the middle of GNNs and their higher-order counterparts. Further, we propose several techniques to aide in choosing the right local graph parameters. Our results connect GNNs with deep results in finite model theory and finite variable logics. Our experimental evaluation shows that adding local graph parameters often has a positive effect for a variety of GNNs, datasets and graph learning tasks.

[Closing the Gap: Tighter Analysis of Alternating Stochastic Gradient Methods for Bilevel Problems](#)

- Tianyi Chen, Yuejiao Sun, Wotao Yin
- abstract@[open-review](#): Stochastic nested optimization, including stochastic compositional, min-max, and bilevel optimization, is gaining popularity in many machine learning applications. While the three problems share a nested structure, existing works often treat them separately, thus developing problem-specific algorithms and analyses. Among various exciting developments, simple SGD-type updates (potentially on multiple variables) are still prevalent in solving this class of nested problems, but they are believed to have a slower convergence rate than non-nested problems. This paper unifies several SGD-type updates for stochastic nested problems into a single SGD approach that we term ALternating Stochastic gradient dEscenT (ALSET) method. By leveraging the hidden smoothness of the problem, this paper presents a tighter analysis of ALSET for stochastic nested problems. Under the new analysis, to achieve an $\$\\epsilon$ -stationary point of the nested problem, it requires $\${\\cal O}(\\epsilon^{-2})$ samples in total. Under certain regularity conditions, applying our results to stochastic compositional, min-max, and reinforcement learning problems either improves or matches the best-known sample complexity in the respective cases. Our results explain why simple SGD-type algorithms in stochastic nested problems all work very well in practice without the need for further modifications.

[Dense Unsupervised Learning for Video Segmentation](#)

- Nikita Araslanov, Simone Schaub-Meyer, Stefan Roth
- abstract@[open-review](#): We present a novel approach to unsupervised learning for video object segmentation (VOS). Unlike previous work, our formulation allows to learn dense feature representations directly in a fully convolutional regime. We rely on uniform grid sampling to extract a set of anchors and train our model to disambiguate between them on both inter- and intra-video levels. However, a naive scheme to train such a model results in a degenerate solution. We propose to prevent this with a simple regularisation scheme, accommodating the equivariance property of the segmentation task to similarity transformations. Our training objective admits efficient implementation and exhibits fast training convergence. On established VOS benchmarks, our approach exceeds the segmentation accuracy of previous work despite using significantly less training data and compute power.

[Charting and Navigating the Space of Solutions for Recurrent Neural Networks](#)

- Elia Turner, Kabir V Dabholkar, Omri Barak
- abstract@[open-review](#): In recent years Recurrent Neural Networks (RNNs) were successfully used to model the way neural activity drives task-related behavior in animals, operating under the implicit assumption that the obtained solutions are universal. Observations in both neuroscience and machine learning challenge this assumption. Animals can approach a given task with a variety of strategies, and training machine learning algorithms introduces the phenomenon of underspecification. These observations imply that every task is associated with a space of solutions. To date, the structure of this space is not understood, limiting the approach of comparing RNNs with neural data. Here, we characterize the space of solutions associated with various tasks. We first study a simple two-neuron network on a task that leads to multiple solutions. We trace the nature of the final solution back to the network’s initial connectivity and identify discrete dynamical regimes that underlie this diversity. We then examine three neuroscience-inspired tasks: Delayed discrimination, Interval discrimination, and Time reproduction. For each task, we find a rich set of solutions. One layer of variability can be found directly in the neural activity of the networks. An additional layer is uncovered by testing the trained networks’ ability to extrapolate, as a perturbation to a system often reveals hidden structure. Furthermore, we relate extrapolation patterns to specific dynamical objects and effective algorithms found by the networks. We introduce a tool to derive the reduced dynamics of networks by generating a compact directed graph describing the essence of the dynamics with regards to behavioral inputs and outputs. Using this representation, we can partition the solutions to each task into a handful of types and show that neural features can partially predict them. Taken together, our results shed light on the concept of the space of solutions and its uses both in Machine learning and in Neuroscience.

[Fast Training Method for Stochastic Compositional Optimization Problems](#)

- Hongchang Gao, Heng Huang
- abstract@[open-review](#): The stochastic compositional optimization problem covers a wide range of machine learning models, such as sparse additive models and model-agnostic meta-learning. Thus, it is necessary to develop efficient methods for its optimization. Existing methods for the stochastic compositional optimization problem only focus on the single machine scenario, which is far from satisfactory when data are distributed on different devices. To address this problem, we propose novel decentralized stochastic compositional gradient descent methods to efficiently train the large-scale stochastic compositional optimization problem. To the best of our knowledge, our work is the first one facilitating decentralized training for this kind of problem. Furthermore, we provide the convergence analysis for our methods, which shows that the convergence rate of our methods can achieve linear speedup with respect to the number of devices. At last, we apply our decentralized training methods to the model-agnostic meta-learning problem, and the experimental results confirm the superior performance of our methods.

[Dual-stream Network for Visual Recognition](#)

- Mingyuan Mao, peng gao, Renrui Zhang, Honghui Zheng, Teli Ma, Yan Peng, Errui Ding, Baochang Zhang, Shumin Han
- abstract@[open-review](#): Transformers with remarkable global representation capacities achieve competitive results for visual tasks, but fail to consider high-level local pattern information in input images. In this paper, we present a generic Dual-stream Network (DS-Net) to fully explore the representation capacity of local and global pattern features for image classification. Our DS-Net can simultaneously calculate fine-grained and integrated features and efficiently fuse them. Specifically, we propose an Intra-scale Propagation module to process two different resolutions in each block and an Inter-Scale Alignment module to perform information interaction across features at dual scales. Besides, we also design a Dual-stream FPN (DS-FPN) to further enhance contextual information for downstream dense predictions. Without bells and whistles, the proposed DS-Net outperforms DeiT-Small by 2.4% in terms of top-1 accuracy on ImageNet-1k and achieves state-of-the-art performance over other Vision Transformers and ResNets. For object detection and instance segmentation, DS-Net-Small respectively outperforms ResNet-50 by 6.4% and 5.5 % in terms of mAP on MSCOCO 2017, and surpasses the previous state-of-the-art scheme, which significantly demonstrates its potential to be a general backbone in vision tasks. The code will be released soon.

[Estimating High Order Gradients of the Data Distribution by Denoising](#)

- Chenlin Meng, Yang Song, Wenzhe Li, Stefano Ermon
- abstract@[open-review](#): The first order derivative of a data density can be estimated efficiently by denoising score matching, and has become an important component in many applications, such as image generation and audio synthesis. Higher order derivatives provide additional local information about the data distribution and enable new applications. Although they can be estimated via automatic differentiation of a learned density model, this can amplify estimation errors and is expensive in high dimensional settings. To overcome these limitations, we propose a method to directly estimate high order derivatives (scores) of a data density from samples. We first show that denoising score matching can be interpreted as a particular case of Tweedie's formula. By leveraging Tweedie's formula on higher order moments, we generalize denoising score matching to estimate higher order derivatives. We demonstrate empirically that models trained with the proposed method can approximate second order derivatives more efficiently and accurately than via automatic differentiation. We show that our models can be used to quantify uncertainty in denoising and to improve the mixing speed of Langevin dynamics via Ozaki discretization for sampling synthetic data and natural images.

[Machine versus Human Attention in Deep Reinforcement Learning Tasks](#)

- Suna (Sihang) Guo, Ruohan Zhang, Bo Liu, Yifeng Zhu, Dana Ballard, Mary Hayhoe, Peter Stone
- abstract@[open-review](#): Deep reinforcement learning (RL) algorithms are powerful tools for solving visuomotor decision tasks. However, the trained models are often difficult to interpret, because they are represented as end-to-end deep neural networks. In this paper, we shed light on the inner workings of such trained models by analyzing the pixels that they attend to during task execution, and comparing them with the pixels attended to by humans executing the same tasks. To this end, we investigate the following two questions that, to the best of our knowledge, have not been previously studied. 1) How similar are the visual representations learned by RL agents and humans when performing the same task? and, 2) How do similarities and differences in these learned representations explain RL agents' performance on these tasks? Specifically, we compare the saliency maps of RL agents against visual attention models of human experts when learning to play Atari games. Further, we analyze how hyperparameters of the deep RL algorithm affect the learned representations and saliency maps of the trained agents. The insights provided have the potential to inform novel algorithms for closing the performance gap between human experts and RL agents.

[Reusing Combinatorial Structure: Faster Iterative Projections over Submodular Base Polytopes](#)

- Jai Moondra, Hassan Mortagy, Swati Gupta
- abstract@[open-review](#): Optimization algorithms such as projected Newton's method, FISTA, mirror descent and its variants enjoy near-optimal regret bounds and convergence rates, but suffer from a computational bottleneck of computing ``projections'' in potentially each iteration (e.g., $\$O(T^{1/2})\$$ regret of online mirror descent). On the other hand, conditional gradient variants solve a linear optimization in each iteration, but result in suboptimal rates (e.g., $\$O(T^{3/4})\$$ regret of online Frank-Wolfe). Motivated by this trade-off in runtime v/s convergence rates, we consider iterative projections of close-by points over widely-prevalent submodular base polytopes $\$B(f)\$$. We develop a toolkit to speed up the computation of projections using both discrete and continuous perspectives. We subsequently adapt the away-step Frank-Wolfe algorithm to use this information and enable early termination. For the special case of cardinality based submodular polytopes, we improve the runtime of computing certain Bregman projections by a factor of $\$O(n/\log(n))\$$. Our theoretical results show orders of magnitude reduction in runtime in preliminary computational experiments.

[Constrained Optimization to Train Neural Networks on Critical and Under-Represented Classes](#)

- Sara Sangalli, Ertunc Erdil, Andeas Hötker, Olivio Donati, Ender Konukoglu
- abstract@[open-review](#): Deep neural networks (DNNs) are notorious for making more mistakes for the classes that have substantially fewer samples than the others during training. Such class imbalance is ubiquitous in clinical applications and very crucial to handle because the classes with fewer samples most often correspond to critical cases (e.g., cancer) where misclassifications can have severe consequences. Not to miss such cases, binary classifiers need to be operated at high True Positive Rates (TPRs) by setting a higher threshold, but this comes at the cost of very high False Positive Rates (FPRs) for problems with class imbalance. Existing methods for learning under class imbalance most often do not take this into account. We argue that prediction accuracy should be improved by emphasizing the reduction of FPRs at high TPRs for problems where misclassification of the positive, i.e. critical, class samples are associated with higher cost. To this end, we pose the training of a DNN for binary classification as a constrained optimization problem and introduce a novel constraint that can be used with existing loss functions to enforce maximal area under the ROC curve (AUC) through prioritizing FPR reduction at high TPR. We solve the resulting constrained optimization problem using an Augmented Lagrangian method (ALM). Going beyond binary, we also propose two possible extensions of the proposed constraint for multi-class classification problems. We present experimental results for image-based binary and multi-class classification applications using an in-house medical imaging dataset, CIFAR10, and CIFAR100. Our results demonstrate that the proposed method improves the baselines in majority of the cases by attaining higher accuracy on critical classes while reducing the misclassification rate for the non-critical class samples.

[Collapsed Variational Bounds for Bayesian Neural Networks](#)

- Marcin Tomczak, Siddharth Swaroop, Andrew Foong, Richard Turner
- abstract@[open-review](#): Recent interest in learning large variational Bayesian Neural Networks (BNNs) has been partly hampered by poor predictive performance caused by underfitting, and their performance is known to be very sensitive to the prior over weights. Current practice often fixes the prior parameters to standard values or tunes them using heuristics or cross-validation. In this paper, we treat prior parameters in a distributional way by extending the model and collapsing the variational bound with respect to their posteriors. This leads to novel and tighter Evidence Lower Bounds (ELBOs) for performing variational inference (VI) in BNNs. Our experiments show that the new bounds significantly improve the performance of Gaussian mean-field VI applied to BNNs on a variety of data sets, demonstrating that mean-field VI works well even in deep models. We also find that the tighter ELBOs can be good optimization targets for learning the hyperparameters of hierarchical priors.

[Consistent Estimation for PCA and Sparse Regression with Oblivious Outliers](#)

- Tommaso d'Orsi, Chih-Hung Liu, Rajai Nasser, Gleb Novikov, David Steurer, Stefan Tiegel
- abstract@[open-review](#): We develop machinery to design efficiently computable and $\$emph{consistent}$ estimators, achieving estimation error approaching zero as the number of observations grows, when facing an oblivious adversary that may corrupt responses in all but an $\$alpha\$$ fraction of the samples. As

concrete examples, we investigate two problems: sparse regression and principal component analysis (PCA). For sparse regression, we achieve consistency for optimal sample size $n \geq \sqrt{k \log d} / \alpha^2$ and optimal error rate $O(\sqrt{(k \log d) / (n \cdot \alpha^2)})$ where n is the number of observations, d is the number of dimensions and k is the sparsity of the parameter vector, allowing the fraction of inliers to be inverse-polynomial in the number of samples. Prior to this work, no estimator was known to be consistent when the fraction of inliers α is $O(1/\log \log n)$, even for (non-spherical) Gaussian design matrices. Results holding under weak design assumptions and in the presence of such general noise have only been shown in dense setting (i.e., general linear regression) very recently by d'Orsi et al.^{[\[ICML-linear-regression\]](#)}. In the context of PCA, we attain optimal error guarantees under broad spikiness assumptions on the parameter matrix (usually used in matrix completion). Previous works could obtain non-trivial guarantees only under the assumptions that the measurement noise corresponding to the inliers is polynomially small in n (e.g., Gaussian with variance $1/n^2$). To devise our estimators, we equip the Huber loss with non-smooth regularizers such as the ℓ_1 norm or the nuclear norm, and extend d'Orsi et al.'s approach^{[\[ICML-linear-regression\]](#)} in a novel way to analyze the loss function. Our machinery appears to be easily applicable to a wide range of estimation problems. We complement these algorithmic results with statistical lower bounds showing that the fraction of inliers that our PCA estimator can deal with is optimal up to a constant factor.

[Offline Constrained Multi-Objective Reinforcement Learning via Pessimistic Dual Value Iteration](#)

- Runzhe Wu, Yufeng Zhang, Zhuoran Yang, Zhaoran Wang
- abstract@[open-review](#): In constrained multi-objective RL, the goal is to learn a policy that achieves the best performance specified by a multi-objective preference function under a constraint. We focus on the offline setting where the RL agent aims to learn the optimal policy from a given dataset. This scenario is common in real-world applications where interactions with the environment are expensive and the constraint violation is dangerous. For such a setting, we transform the original constrained problem into a primal-dual formulation, which is solved via dual gradient ascent. Moreover, we propose to combine such an approach with pessimism to overcome the uncertainty in offline data, which leads to our Pessimistic Dual Iteration (PEDI). We establish upper bounds on both the suboptimality and constraint violation for the policy learned by PEDI based on an arbitrary dataset, which proves that PEDI is provably sample efficient. We also specialize PEDI to the setting with linear function approximation. To the best of our knowledge, we propose the first provably efficient constrained multi-objective RL algorithm with offline data without any assumption on the coverage of the dataset.

[Absolute Neighbour Difference based Correlation Test for Detecting Heteroscedastic Relationships](#)

- Lifeng Zhang
- abstract@[open-review](#): It is a challenge to detect complicated data relationships thoroughly. Here, we propose a new statistical measure, named the absolute neighbour difference based neighbour correlation coefficient, to detect the associations between variables through examining the heteroscedasticity of the unpredictable variation of dependent variables. Different from previous studies, the new method concentrates on measuring nonfunctional relationships rather than functional or mixed associations. Either used alone or in combination with other measures, it enables not only a convenient test of heteroscedasticity, but also measuring functional and nonfunctional relationships separately that obviously leads to a deeper insight into the data associations. The method is concise and easy to implement that does not rely on explicitly estimating the regression residuals or the dependencies between variables so that it is not restricted to any kind of model assumption. The mechanisms of the correlation test are proved in theory and demonstrated with numerical analyses.

[Batch Multi-Fidelity Bayesian Optimization with Deep Auto-Regressive Networks](#)

- Shibo Li, Robert Kirby, Shandian Zhe
- abstract@[open-review](#): Bayesian optimization (BO) is a powerful approach for optimizing black-box, expensive-to-evaluate functions. To enable a flexible trade-off between the cost and accuracy, many applications allow the function to be evaluated at different fidelities. In order to reduce the optimization cost while maximizing the benefit-cost ratio, in this paper we propose Batch Multi-fidelity Bayesian Optimization with Deep Auto-Regressive Networks (BMBO-DARN). We use a set of Bayesian neural networks to construct a fully auto-regressive model, which is expressive enough to capture strong yet complex relationships across all the fidelities, so as to improve the surrogate learning and optimization performance. Furthermore, to enhance the quality and diversity of queries, we develop a simple yet efficient batch querying method, without any combinatorial search over the fidelities. We propose a batch acquisition function based on Max-value Entropy Search (MES) principle, which penalizes highly correlated queries and encourages diversity. We use posterior samples and moment matching to fulfill efficient computation of the acquisition function, and conduct alternating optimization over every fidelity-input pair, which guarantees an improvement at each step. We demonstrate the advantage of our approach on four real-world hyperparameter optimization applications.

[Mastering Atari Games with Limited Data](#)

- Weirui Ye, Shaohuai Liu, Thanard Kurutach, Pieter Abbeel, Yang Gao
- abstract@[open-review](#): Reinforcement learning has achieved great success in many applications. However, sample efficiency remains a key challenge, with prominent methods requiring millions (or even billions) of environment steps to train. Recently, there has been significant progress in sample efficient image-based RL algorithms; however, consistent human-level performance on the Atari game benchmark remains an elusive goal. We propose a sample efficient model-based visual RL algorithm built on MuZero, which we name EfficientZero. Our method achieves 194.3% mean human performance and 109.0% median performance on the Atari 100k benchmark with only two hours of real-time game experience and outperforms the state-of-the-art SAC in some tasks on the DMControl 100k benchmark. This is the first time an algorithm achieves super-human performance on Atari games with such little data. EfficientZero's performance is also close to DQN's performance at 200 million frames while we consume 500 times less data. EfficientZero's low sample complexity and high performance can bring RL closer to real-world applicability. We implement our algorithm in an easy-to-understand manner and it is available at <https://github.com/YeWR/EfficientZero>. We hope it will accelerate the research of MCTS-based RL algorithms in the wider community.

[Dealing With Misspecification In Fixed-Confidence Linear Top-m Identification](#)

- Clémence Réda, Andrea Tirinzoni, Rémy Degenne
- abstract@[open-review](#): We study the problem of the identification of m arms with largest means under a fixed error rate δ (fixed-confidence Top- m identification), for misspecified linear bandit models. This problem is motivated by practical applications, especially in medicine and recommendation systems, where linear models are popular due to their simplicity and the existence of efficient algorithms, but in which data inevitably deviates from linearity. In this work, we first derive a tractable lower bound on the sample complexity of any δ -correct algorithm for the general Top- m identification problem. We show that knowing the scale of the deviation from linearity is necessary to exploit the structure of the problem. We then describe the first algorithm for this setting, which is both practical and adapts to the amount of misspecification. We derive an upper bound to its sample complexity which confirms this adaptivity and that matches the lower bound when $\delta \rightarrow 0$. Finally, we evaluate our algorithm on both synthetic and real-world data, showing competitive performance with respect to existing baselines.

[Why Generalization in RL is Difficult: Epistemic POMDPs and Implicit Partial Observability](#)

- Dibya Ghosh, Jad Rahme, Aviral Kumar, Amy Zhang, Ryan P. Adams, Sergey Levine
- abstract@[open-review](#): Generalization is a central challenge for the deployment of reinforcement learning (RL) systems in the real world. In this paper, we show that the sequential structure of the RL problem necessitates new approaches to generalization beyond the well-studied techniques used in supervised

learning. While supervised learning methods can generalize effectively without explicitly accounting for epistemic uncertainty, we describe why appropriate uncertainty handling can actually be essential in RL. We show that generalization to unseen test conditions from a limited number of training conditions induces a kind of implicit partial observability, effectively turning even fully-observed MDPs into POMDPs. Informed by this observation, we recast the problem of generalization in RL as solving the induced partially observed Markov decision process, which we call the epistemic POMDP. We demonstrate the failure modes of algorithms that do not appropriately handle this partial observability, and suggest a simple ensemble-based technique for approximately solving the partially observed problem. Empirically, we demonstrate that our simple algorithm derived from the epistemic POMDP achieves significant gains in generalization over current methods on the Procgen benchmark suite.

[Set Prediction in the Latent Space](#)

- Konpat Preechakul, Chawan Piansaddhayanon, Burin Naowarat, Tirasan Khandhawit, Sira Sriswasdi, Ekapol Chuangsuwanich
- abstract@[open-review](#): Set prediction tasks require the matching between predicted set and ground truth set in order to propagate the gradient signal. Recent works have performed this matching in the original feature space thus requiring predefined distance functions. We propose a method for learning the distance function by performing the matching in the latent space learned from encoding networks. This method enables the use of teacher forcing which was not possible previously since matching in the feature space must be computed after the entire output sequence is generated. Nonetheless, a naive implementation of latent set prediction might not converge due to permutation instability. To address this problem, we provide sufficient conditions for permutation stability which begets an algorithm to improve the overall model convergence. Experiments on several set prediction tasks, including image captioning and object detection, demonstrate the effectiveness of our method.

[Best of Both Worlds: Practical and Theoretically Optimal Submodular Maximization in Parallel](#)

- Yixin Chen, Tonmoy Dey, Alan Kuhnle
- abstract@[open-review](#): For the problem of maximizing a monotone, submodular function with respect to a cardinality constraint $\$k\$$ on a ground set of size $\$n\$$, we provide an algorithm that achieves the state-of-the-art in both its empirical performance and its theoretical properties, in terms of adaptive complexity, query complexity, and approximation ratio; that is, it obtains, with high probability, query complexity of $\$O(n)\$$ in expectation, adaptivity of $\$O(\log(n))\$$, and approximation ratio of nearly $\$1-1/e\$$. The main algorithm is assembled from two components which may be of independent interest. The first component of our algorithm, LINEARSEQ, is useful as a preprocessing algorithm to improve the query complexity of many algorithms. Moreover, a variant of LINEARSEQ is shown to have adaptive complexity of $\$O(\sqrt{\log(n/k)})\$$ which is smaller than that of any previous algorithm in the literature. The second component is a parallelizable thresholding procedure THRESHOLDSEQ for adding elements with gain above a constant threshold. Finally, we demonstrate that our main algorithm empirically outperforms, in terms of runtime, adaptive rounds, total queries, and objective values, the previous state-of-the-art algorithm FAST in a comprehensive evaluation with six submodular objective functions.

[Fine-grained Generalization Analysis of Inductive Matrix Completion](#)

- Antoine Ledent, Rodrigo Alves, Yunwen Lei, Marius Kloft
- abstract@[open-review](#): In this paper, we bridge the gap between the state-of-the-art theoretical results for matrix completion with the nuclear norm and their equivalent in \textit{inductive matrix completion}: (1) In the distribution-free setting, we prove bounds improving the previously best scaling of $\$O(r d^2)\$$ to $\$O(d^{3/2} \sqrt{r})\$$, where $\$d\$$ is the dimension of the side information and $\$r\$$ is the rank. (2) We introduce the (smoothed) \textit{adjusted trace-norm minimization} strategy, an inductive analogue of the weighted trace norm, for which we show guarantees of the order $\$O(dr)\$$ under arbitrary sampling. In the inductive case, a similar rate was previously achieved only under uniform sampling and for exact recovery. Both our results align with the state of the art in the particular case of standard (non-inductive) matrix completion, where they are known to be tight up to log terms. Experiments further confirm that our strategy outperforms standard inductive matrix completion on various synthetic datasets and real problems, justifying its place as an important tool in the arsenal of methods for matrix completion using side information.

[Learning Frequency Domain Approximation for Binary Neural Networks](#)

- Yixing Xu, Kai Han, Chang Xu, Yehui Tang, Chunjing XU, Yunhe Wang
- abstract@[open-review](#): Binary neural networks (BNNs) represent original full-precision weights and activations into 1-bit with sign function. Since the gradient of the conventional sign function is almost zero everywhere which cannot be used for back-propagation, several attempts have been proposed to alleviate the optimization difficulty by using approximate gradient. However, those approximations corrupt the main direction of factual gradient. To this end, we propose to estimate the gradient of sign function in the Fourier frequency domain using the combination of sine functions for training BNNs, namely frequency domain approximation (FDA). The proposed approach does not affect the low-frequency information of the original sign function which occupies most of the overall energy, and high-frequency coefficients will be ignored to avoid the huge computational overhead. In addition, we embed a noise adaptation module into the training phase to compensate the approximation error. The experiments on several benchmark datasets and neural architectures illustrate that the binary network learned using our method achieves the state-of-the-art accuracy. Code will be available at <https://gitee.com/mindspore/models/tree/master/research/cv/FDA-BNN>.

[Reformulating Zero-shot Action Recognition for Multi-label Actions](#)

- Alec Kerrigan, Kevin Duarte, Yogesh Rawat, Mubarak Shah
- abstract@[open-review](#): The goal of zero-shot action recognition (ZSAR) is to classify action classes which were not previously seen during training. Traditionally, this is achieved by training a network to map, or regress, visual inputs to a semantic space where a nearest neighbor classifier is used to select the closest target class. We argue that this approach is sub-optimal due to the use of nearest neighbor on static semantic space and is ineffective when faced with multi-label videos - where two semantically distinct co-occurring action categories cannot be predicted with high confidence. To overcome these limitations, we propose a ZSAR framework which does not rely on nearest neighbor classification, but rather consists of a pairwise scoring function. Given a video and a set of action classes, our method predicts a set of confidence scores for each class independently. This allows for the prediction of several semantically distinct classes within one video input. Our evaluations show that our method not only achieves strong performance on three single-label action classification datasets (UCF-101, HMDB, and RareAct), but also outperforms previous ZSAR approaches on a challenging multi-label dataset (AVA) and a real-world surprise activity detection dataset (MEVA).

[Optimal Best-Arm Identification Methods for Tail-Risk Measures](#)

- Shubhada Agrawal, Wouter M. Koolen, Sandeep Juneja
- abstract@[open-review](#): Conditional value-at-risk (CVaR) and value-at-risk (VaR) are popular tail-risk measures in finance and insurance industries as well as in highly reliable, safety-critical uncertain environments where often the underlying probability distributions are heavy-tailed. We use the multi-armed bandit best-arm identification framework and consider the problem of identifying the arm from amongst finitely many that has the smallest CVaR, VaR, or weighted sum of CVaR and mean. The latter captures the risk-return trade-off common in finance. Our main contribution is an optimal $\$delta\$$ -correct algorithm that acts on general arms, including heavy-tailed distributions, and matches the lower bound on the expected number of samples needed, asymptotically (as $\$delta\$$ approaches $\$0\$$). The algorithm requires solving a non-convex optimization problem in the space of probability measures, that requires delicate analysis. En-route, we develop new non-asymptotic, anytime-valid, empirical-likelihood-based concentration inequalities for tail-risk measures.

SyMetric: Measuring the Quality of Learnt Hamiltonian Dynamics Inferred from Vision

- Irina Higgins, Peter Wirsberger, Andrew Jaegle, Aleksandar Botev
- abstract@[open-review](#): A recently proposed class of models attempts to learn latent dynamics from high-dimensional observations, like images, using priors informed by Hamiltonian mechanics. While these models have important potential applications in areas like robotics or autonomous driving, there is currently no good way to evaluate their performance: existing methods primarily rely on image reconstruction quality, which does not always reflect the quality of the learnt latent dynamics. In this work, we empirically highlight the problems with the existing measures and develop a set of new measures, including a binary indicator of whether the underlying Hamiltonian dynamics have been faithfully captured, which we call Symplecticity Metric or SyMetric. Our measures take advantage of the known properties of Hamiltonian dynamics and are more discriminative of the model's ability to capture the underlying dynamics than reconstruction error. Using SyMetric, we identify a set of architectural choices that significantly improve the performance of a previously proposed model for inferring latent dynamics from pixels, the Hamiltonian Generative Network (HGN). Unlike the original HGN, the new SyMetric is able to discover an interpretable phase space with physically meaningful latents on some datasets. Furthermore, it is stable for significantly longer rollouts on a diverse range of 13 datasets, producing rollouts of essentially infinite length both forward and backwards in time with no degradation in quality on a subset of the datasets.

Learning with Holographic Reduced Representations

- Ashwinkumar Ganesan, Hang Gao, Sunil Gandhi, Edward Raff, Tim Oates, James Holt, Mark McLean
- abstract@[open-review](#): Holographic Reduced Representations (HRR) are a method for performing symbolic AI on top of real-valued vectors by associating each vector with an abstract concept, and providing mathematical operations to manipulate vectors as if they were classic symbolic objects. This method has seen little use outside of older symbolic AI work and cognitive science. Our goal is to revisit this approach to understand if it is viable for enabling a hybrid neural-symbolic approach to learning as a differential component of a deep learning architecture. HRRs today are not effective in a differential solution due to numerical instability, a problem we solve by introducing a projection step that forces the vectors to exist in a well behaved point in space. In doing so we improve the concept retrieval efficacy of HRRs by over \$100\times\$. Using multi-label classification we demonstrate how to leverage the symbolic HRR properties to develop a output layer and loss function that is able to learn effectively, and allows us to investigate some of the pros and cons of an HRR neuro-symbolic learning approach.

Learning Barrier Certificates: Towards Safe Reinforcement Learning with Zero Training-time Violations

- Yuping Luo, Tengyu Ma
- abstract@[open-review](#): Training-time safety violations have been a major concern when we deploy reinforcement learning algorithms in the real world. This paper explores the possibility of safe RL algorithms with zero training-time safety violations in the challenging setting where we are only given a safe but trivial-reward initial policy without any prior knowledge of the dynamics and additional offline data. We propose an algorithm, Co-trained Barrier Certificate for Safe RL (CRABS), which iteratively learns barrier certificates, dynamics models, and policies. The barrier certificates are learned via adversarial training and ensure the policy's safety assuming calibrated learned dynamics. We also add a regularization term to encourage larger certified regions to enable better exploration. Empirical simulations show that zero safety violations are already challenging for a suite of simple environments with only 2-4 dimensional state space, especially if high-reward policies have to visit regions near the safety boundary. Prior methods require hundreds of violations to achieve decent rewards on these tasks, whereas our proposed algorithms incur zero violations.

On the Second-order Convergence Properties of Random Search Methods

- Aurelien Lucchi, Antonio Orvieto, Adamos Solomou
- abstract@[open-review](#): We study the theoretical convergence properties of random-search methods when optimizing non-convex objective functions without having access to derivatives. We prove that standard random-search methods that do not rely on second-order information converge to a second-order stationary point. However, they suffer from an exponential complexity in terms of the input dimension of the problem. In order to address this issue, we propose a novel variant of random search that exploits negative curvature by only relying on function evaluations. We prove that this approach converges to a second-order stationary point at a much faster rate than vanilla methods: namely, the complexity in terms of the number of function evaluations is only linear in the problem dimension. We test our algorithm empirically and find good agreements with our theoretical results.

Noether's Learning Dynamics: Role of Symmetry Breaking in Neural Networks

- Hidenori Tanaka, Daniel Kunin
- abstract@[open-review](#): In nature, symmetry governs regularities, while symmetry breaking brings texture. In artificial neural networks, symmetry has been a central design principle to efficiently capture regularities in the world, but the role of symmetry breaking is not well understood. Here, we develop a theoretical framework to study the "geometry of learning dynamics" in neural networks, and reveal a key mechanism of explicit symmetry breaking behind the efficiency and stability of modern neural networks. To build this understanding, we model the discrete learning dynamics of gradient descent using a continuous-time Lagrangian formulation, in which the learning rule corresponds to the kinetic energy and the loss function corresponds to the potential energy. Then, we identify "kinetic symmetry breaking" (KSB), the condition when the kinetic energy explicitly breaks the symmetry of the potential function. We generalize Noether's theorem known in physics to take into account KSB and derive the resulting motion of the Noether charge: "Noether's Learning Dynamics" (NLD). Finally, we apply NLD to neural networks with normalization layers and reveal how KSB introduces a mechanism of implicit adaptive optimization, establishing an analogy between learning dynamics induced by normalization layers and RMSProp. Overall, through the lens of Lagrangian mechanics, we have established a theoretical foundation to discover geometric design principles for the learning dynamics of neural networks.

A Theory of the Distortion-Perception Tradeoff in Wasserstein Space

- Dror Freirich, Tomer Michaeli, Ron Meir
- abstract@[open-review](#): The lower the distortion of an estimator, the more the distribution of its outputs generally deviates from the distribution of the signals it attempts to estimate. This phenomenon, known as the perception-distortion tradeoff, has captured significant attention in image restoration, where it implies that fidelity to ground truth images comes on the expense of perceptual quality (deviation from statistics of natural images). However, despite the increasing popularity of performing comparisons on the perception-distortion plane, there remains an important open question: what is the minimal distortion that can be achieved under a given perception constraint? In this paper, we derive a closed form expression for this distortion-perception (DP) function for the mean squared-error (MSE) distortion and Wasserstein-2 perception index. We prove that the DP function is always quadratic, regardless of the underlying distribution. This stems from the fact that estimators on the DP curve form a geodesic in Wasserstein space. In the Gaussian setting, we further provide a closed form expression for such estimators. For general distributions, we show how these estimators can be constructed from the estimators at the two extremes of the tradeoff: The global MSE minimizer, and a minimizer of the MSE under a perfect perceptual quality constraint. The latter can be obtained as a stochastic transformation of the former.

Neural Production Systems

- Anirudh Goyal ALIAS PARTH GOYAL, Aniket Didolkar, Nan Rosemary Ke, Charles Blundell, Philippe Beaudoin, Nicolas Heess, Michael C. Mozer, Yoshua Bengio
- abstract@[open-review](#): Visual environments are structured, consisting of distinct objects or entities. These entities have properties---visible or latent---that determine the manner in which they interact with one another. To partition images into entities, deep-learning researchers have proposed structural inductive biases such as slot-based architectures. To model interactions among entities, equivariant graph neural nets (GNNs) are used, but these are not particularly well suited to the task for two reasons. First, GNNs do not predispose interactions to be sparse, as relationships among independent entities are likely to be. Second, GNNs do not factorize knowledge about interactions in an entity-conditional manner. As an alternative, we take inspiration from cognitive science and resurrect a classic approach, production systems, which consist of a set of rule templates that are applied by binding placeholder variables in the rules to specific entities. Rules are scored on their match to entities, and the best fitting rules are applied to update entity properties. In a series of experiments, we demonstrate that this architecture achieves a flexible, dynamic flow of control and serves to factorize entity-specific and rule-based information. This disentangling of knowledge achieves robust future-state prediction in rich visual environments, outperforming state-of-the-art methods using GNNs, and allows for the extrapolation from simple (few object) environments to more complex environments.

[Smoothness Matrices Beat Smoothness Constants: Better Communication Compression Techniques for Distributed Optimization](#)

- Mher Safaryan, Filip Hanzely, Peter Richtarik
- abstract@[open-review](#): Large scale distributed optimization has become the default tool for the training of supervised machine learning models with a large number of parameters and training data. Recent advancements in the field provide several mechanisms for speeding up the training, including {\em compressed communication}, {\em variance reduction} and {\em acceleration}. However, none of these methods is capable of exploiting the inherently rich data-dependent smoothness structure of the local losses beyond standard smoothness constants. In this paper, we argue that when training supervised models, {\em smoothness matrices}---information-rich generalizations of the ubiquitous smoothness constants---can and should be exploited for further dramatic gains, both in theory and practice. In order to further alleviate the communication burden inherent in distributed optimization, we propose a novel communication sparsification strategy that can take full advantage of the smoothness matrices associated with local losses. To showcase the power of this tool, we describe how our sparsification technique can be adapted to three distributed optimization algorithms---DCGD, DIANA and ADIANA---yielding significant savings in terms of communication complexity. The new methods always outperform the baselines, often dramatically so.

[Increasing Liquid State Machine Performance with Edge-of-Chaos Dynamics Organized by Astrocyte-modulated Plasticity](#)

- Vladimir Ivanov, Konstantinos Michmizos
- abstract@[open-review](#): The liquid state machine (LSM) combines low training complexity and biological plausibility, which has made it an attractive machine learning framework for edge and neuromorphic computing paradigms. Originally proposed as a model of brain computation, the LSM tunes its internal weights without backpropagation of gradients, which results in lower performance compared to multi-layer neural networks. Recent findings in neuroscience suggest that astrocytes, a long-neglected non-neuronal brain cell, modulate synaptic plasticity and brain dynamics, tuning brain networks to the vicinity of the computationally optimal critical phase transition between order and chaos. Inspired by this disruptive understanding of how brain networks self-tune, we propose the neuron-astrocyte liquid state machine (NALSM) that addresses under-performance through self-organized near-critical dynamics. Similar to its biological counterpart, the astrocyte model integrates neuronal activity and provides global feedback to spike-timing-dependent plasticity (STDP), which self-organizes NALSM dynamics around a critical branching factor that is associated with the edge-of-chaos. We demonstrate that NALSM achieves state-of-the-art accuracy versus comparable LSM methods, without the need for data-specific hand-tuning. With a top accuracy of \$97.61\%\$ on MNIST, \$97.51\%\$ on N-MNIST, and \$85.84\%\$ on Fashion-MNIST, NALSM achieved comparable performance to current fully-connected multi-layer spiking neural networks trained via backpropagation. Our findings suggest that the further development of brain-inspired machine learning methods has the potential to reach the performance of deep learning, with the added benefits of supporting robust and energy-efficient neuromorphic computing on the edge.

[Fair Sortition Made Transparent](#)

- Bailey Flanigan, Gregory Kehne, Ariel D. Procaccia
- abstract@[open-review](#): Sortition is an age-old democratic paradigm, widely manifested today through the random selection of citizens' assemblies. Recently-deployed algorithms select assemblies \textit{maximally fairly}, meaning that subject to demographic quotas, they give all potential participants as equal a chance as possible of being chosen. While these fairness gains can bolster the legitimacy of citizens' assemblies and facilitate their uptake, existing algorithms remain limited by their lack of transparency. To overcome this hurdle, in this work we focus on panel selection by uniform lottery, which is easy to realize in an observable way. By this approach, the final assembly is selected by uniformly sampling some pre-selected set of \$m\$ possible assemblies. We provide theoretical guarantees on the fairness attainable via this type of uniform lottery, as compared to the existing maximally fair but opaque algorithms, for two different fairness objectives. We complement these results with experiments on real-world instances that demonstrate the viability of the uniform lottery approach as a method of selecting assemblies both fairly and transparently.

[A Max-Min Entropy Framework for Reinforcement Learning](#)

- Seungyul Han, Youngchul Sung
- abstract@[open-review](#): In this paper, we propose a max-min entropy framework for reinforcement learning (RL) to overcome the limitation of the soft actor-critic (SAC) algorithm implementing the maximum entropy RL in model-free sample-based learning. Whereas the maximum entropy RL guides learning for policies to reach states with high entropy in the future, the proposed max-min entropy framework aims to learn to visit states with low entropy and maximize the entropy of these low-entropy states to promote better exploration. For general Markov decision processes (MDPs), an efficient algorithm is constructed under the proposed max-min entropy framework based on disentanglement of exploration and exploitation. Numerical results show that the proposed algorithm yields drastic performance improvement over the current state-of-the-art RL algorithms.

[Reward is enough for convex MDPs](#)

- Tom Zahavy, Brendan O'Donoghue, Guillaume Desjardins, Satinder Singh
- abstract@[open-review](#): Maximising a cumulative reward function that is Markov and stationary, i.e., defined over state-action pairs and independent of time, is sufficient to capture many kinds of goals in a Markov decision process (MDP). However, not all goals can be captured in this manner. In this paper we study convex MDPs in which goals are expressed as convex functions of the stationary distribution and show that they cannot be formulated using stationary reward functions. Convex MDPs generalize the standard reinforcement learning (RL) problem formulation to a larger framework that includes many supervised and unsupervised RL problems, such as apprenticeship learning, constrained MDPs, and so-called pure exploration'. Our approach is to reformulate the convex MDP problem as a min-max game involving policy and cost (negative reward) players', using Fenchel duality. We propose a meta-algorithm for solving this problem and show that it unifies many existing algorithms in the literature.

[Fast Doubly-Adaptive MCMC to Estimate the Gibbs Partition Function with Weak Mixing Time Bounds](#)

- Shahrzad Haddadan, Yue Zhuang, Cyrus Cousins, Eli Upfal
- abstract@[open-review](#): We present a novel method for reducing the computational complexity of rigorously estimating the partition functions of Gibbs (or Boltzmann) distributions, which arise ubiquitously in probabilistic graphical models. A major obstacle to applying the Gibbs distribution in practice is the need to estimate their partition function (normalizing constant). The state of the art in addressing this problem is multi-stage algorithms which consist of a

cooling schedule and a mean estimator in each step of the schedule. While the cooling schedule in these algorithms is adaptive, the mean estimate computations use MCMC as a black-box to draw approximately-independent samples. Here we develop a doubly adaptive approach, combining the adaptive cooling schedule with an adaptive MCMC mean estimator, whose number of Markov chain steps adapts dynamically to the underlying chain. Through rigorous theoretical analysis, we prove that our method outperforms the state of the art algorithms in several factors: (1) The computational complexity of our method is smaller; (2) Our method is less sensitive to loose bounds on mixing times, an inherent components in these algorithms; and (3) The improvement obtained by our method is particularly significant in the most challenging regime of high precision estimates. We demonstrate the advantage of our method in experiments run on classic factor graphs, such as voting models and Ising models.

[Does enforcing fairness mitigate biases caused by subpopulation shift?](#)

- Subha Maity, Debarghya Mukherjee, Mikhail Yurochkin, Yuekai Sun
- abstract@[open-review](#): Many instances of algorithmic bias are caused by subpopulation shifts. For example, ML models often perform worse on demographic groups that are underrepresented in the training data. In this paper, we study whether enforcing algorithmic fairness during training improves the performance of the trained model in the \emph{target domain}. On one hand, we conceive scenarios in which enforcing fairness does not improve performance in the target domain. In fact, it may even harm performance. On the other hand, we derive necessary and sufficient conditions under which enforcing algorithmic fairness leads to the Bayes model in the target domain. We also illustrate the practical implications of our theoretical results in simulations and on real data.

[Implicit Deep Adaptive Design: Policy-Based Experimental Design without Likelihoods](#)

- Desi R Ivanova, Adam Foster, Steven Kleinegesse, Michael U. Gutmann, Thomas Rainforth
- abstract@[open-review](#): We introduce implicit Deep Adaptive Design (iDAD), a new method for performing adaptive experiments in real-time with implicit models. iDAD amortizes the cost of Bayesian optimal experimental design (BOED) by learning a design policy network upfront, which can then be deployed quickly at the time of the experiment. The iDAD network can be trained on any model which simulates differentiable samples, unlike previous design policy work that requires a closed form likelihood and conditionally independent experiments. At deployment, iDAD allows design decisions to be made in milliseconds, in contrast to traditional BOED approaches that require heavy computation during the experiment itself. We illustrate the applicability of iDAD on a number of experiments, and show that it provides a fast and effective mechanism for performing adaptive design with implicit models.

[Sample-Efficient Learning of Stackelberg Equilibria in General-Sum Games](#)

- Yu Bai, Chi Jin, Huan Wang, Caiming Xiong
- abstract@[open-review](#): Real world applications such as economics and policy making often involve solving multi-agent games with two unique features: (1) The agents are inherently asymmetric and partitioned into leaders and followers; (2) The agents have different reward functions, thus the game is general-sum. The majority of existing results in this field focuses on either symmetric solution concepts (e.g. Nash equilibrium) or zero-sum games. It remains open how to learn the Stackelberg equilibrium---an asymmetric analog of the Nash equilibrium---in general-sum games efficiently from noisy samples. This paper initiates the theoretical study of sample-efficient learning of the Stackelberg equilibrium, in the bandit feedback setting where we only observe noisy samples of the reward. We consider three representative two-player general-sum games: bandit games, bandit-reinforcement learning (bandit-RL) games, and linear bandit games. In all these games, we identify a fundamental gap between the exact value of the Stackelberg equilibrium and its estimated version using finitely many noisy samples, which can not be closed information-theoretically regardless of the algorithm. We then establish sharp positive results on sample-efficient learning of Stackelberg equilibrium with value optimal up to the gap identified above, with matching lower bounds in the dependency on the gap, error tolerance, and the size of the action spaces. Overall, our results unveil unique challenges in learning Stackelberg equilibria under noisy bandit feedback, which we hope could shed light on future research on this topic.

[Non-approximate Inference for Collective Graphical Models on Path Graphs via Discrete Difference of Convex Algorithm](#)

- Yasunori Akagi, Naoki Marumo, Hideaki Kim, Takeshi Kurashima, Hiroyuki Toda
- abstract@[open-review](#): The importance of aggregated count data, which is calculated from the data of multiple individuals, continues to increase. Collective Graphical Model (CGM) is a probabilistic approach to the analysis of aggregated data. One of the most important operations in CGM is maximum a posteriori (MAP) inference of unobserved variables under given observations. Because the MAP inference problem for general CGMs has been shown to be NP-hard, an approach that solves an approximate problem has been proposed. However, this approach has two major drawbacks. First, the quality of the solution deteriorates when the values in the count tables are small, because the approximation becomes inaccurate. Second, since continuous relaxation is applied, the integrality constraints of the output are violated. To resolve these problems, this paper proposes a new method for MAP inference for CGMs on path graphs. Our method is based on the Difference of Convex Algorithm (DCA), which is a general methodology to minimize a function represented as the sum of a convex function and a concave function. In our algorithm, important subroutines in DCA can be efficiently calculated by minimum convex cost flow algorithms. Experiments show that the proposed method outputs higher quality solutions than the conventional approach.

[Implicit Task-Driven Probability Discrepancy Measure for Unsupervised Domain Adaptation](#)

- Mao Li, Kaiqi Jiang, Xinhua Zhang
- abstract@[open-review](#): Probability discrepancy measure is a fundamental construct for numerous machine learning models such as weakly supervised learning and generative modeling. However, most measures overlook the fact that the distributions are not the end-product of learning, but are the basis of downstream predictor. Therefore it is important to warp the probability discrepancy measure towards the end tasks, and we hence propose a new bi-level optimization based approach so that the two distributions are compared not uniformly against the entire hypothesis space, but only with respect to the optimal predictor for the downstream end task. When applied to margin disparity discrepancy and contrastive domain discrepancy, our method significantly improves the performance in unsupervised domain adaptation, and enjoys a much more principled training process.

[SBO-RNN: Reformulating Recurrent Neural Networks via Stochastic Bilevel Optimization](#)

- Ziming Zhang, Yun Yue, Guojun Wu, Yanhua Li, Haichong Zhang
- abstract@[open-review](#): In this paper we consider the training stability of recurrent neural networks (RNNs) and propose a family of RNNs, namely SBO-RNN, that can be formulated using stochastic bilevel optimization (SBO). With the help of stochastic gradient descent (SGD), we manage to convert the SBO problem into an RNN where the feedforward and backpropagation solve the lower and upper-level optimization for learning hidden states and their hyperparameters, respectively. We prove that under mild conditions there is no vanishing or exploding gradient in training SBO-RNN. Empirically we demonstrate our approach with superior performance on several benchmark datasets, with fewer parameters, less training data, and much faster convergence. Code is available at <https://zhang-vislab.github.io>.

[Navigating to the Best Policy in Markov Decision Processes](#)

- Aymen Al Marjani, Aurélien Garivier, Alexandre Proutiere

- abstract@[open-review](#): We investigate the classical active pure exploration problem in Markov Decision Processes, where the agent sequentially selects actions and, from the resulting system trajectory, aims at identifying the best policy as fast as possible. We propose a problem-dependent lower bound on the average number of steps required before a correct answer can be given with probability at least $1 - \delta$. We further provide the first algorithm with an instance-specific sample complexity in this setting. This algorithm addresses the general case of communicating MDPs; we also propose a variant with a reduced exploration rate (and hence faster convergence) under an additional ergodicity assumption. This work extends previous results relative to the generative setting~\cite{pmlr-v139-marjani21a}, where the agent could at each step query the random outcome of any (state, action) pair. In contrast, we show here how to deal with the navigation constraints, induced by the online setting. Our analysis relies on an ergodic theorem for non-homogeneous Markov chains which we consider of wide interest in the analysis of Markov Decision Processes.

[A Faster Decentralized Algorithm for Nonconvex Minimax Problems](#)

- Wenhan Xian, Feihu Huang, Yanfu Zhang, Heng Huang
- abstract@[open-review](#): In this paper, we study the nonconvex-strongly-concave minimax optimization problem on decentralized setting. The minimax problems are attracting increasing attentions because of their popular practical applications such as policy evaluation and adversarial training. As training data become larger, distributed training has been broadly adopted in machine learning tasks. Recent research works show that the decentralized distributed data-parallel training techniques are specially promising, because they can achieve the efficient communications and avoid the bottleneck problem on the central node or the latency of low bandwidth network. However, the decentralized minimax problems were seldom studied in literature and the existing methods suffer from very high gradient complexity. To address this challenge, we propose a new faster decentralized algorithm, named as DM-HSGD, for nonconvex minimax problems by using the variance reduced technique of hybrid stochastic gradient descent. We prove that our DM-HSGD algorithm achieves stochastic first-order oracle (SFO) complexity of $O(\kappa^3 \epsilon^{-3})$ for decentralized stochastic nonconvex-strongly-concave problem to search an ϵ -stationary point, which improves the exiting best theoretical results. Moreover, we also prove that our algorithm achieves linear speedup with respect to the number of workers. Our experiments on decentralized settings show the superior performance of our new algorithm.

[Generalization Bounds For Meta-Learning: An Information-Theoretic Analysis](#)

- Qi CHEN, Changjian Shui, Mario Marchand
- abstract@[open-review](#): We derive a novel information-theoretic analysis of the generalization property of meta-learning algorithms. Concretely, our analysis proposes a generic understanding in both the conventional learning-to-learn framework \cite{amit2018meta} and the modern model-agnostic meta-learning (MAML) algorithms \cite{finn2017model}. Moreover, we provide a data-dependent generalization bound for the stochastic variant of MAML, which is non-vacuous for deep few-shot learning. As compared to previous bounds that depend on the square norms of gradients, empirical validations on both simulated data and a well-known few-shot benchmark show that our bound is orders of magnitude tighter in most conditions.

[ReLU Regression with Massart Noise](#)

- Ilias Diakonikolas, Jong Ho Park, Christos Tzamos
- abstract@[open-review](#): We study the fundamental problem of ReLU regression, where the goal is to fit Rectified Linear Units (ReLUs) to data. This supervised learning task is efficiently solvable in the realizable setting, but is known to be computationally hard with adversarial label noise. In this work, we focus on ReLU regression in the Massart noise model, a natural and well-studied semi-random noise model. In this model, the label of every point is generated according to a function in the class, but an adversary is allowed to change this value arbitrarily with some probability, which is at most $\eta < 1/2$. We develop an efficient algorithm that achieves exact parameter recovery in this model under mild anti-concentration assumptions on the underlying distribution. Such assumptions are necessary for exact recovery to be information-theoretically possible. We demonstrate that our algorithm significantly outperforms naive applications of ℓ_1 and ℓ_2 regression on both synthetic and real data.

[Identification of the Generalized Condorcet Winner in Multi-dueling Bandits](#)

- Björn Haddenhurst, Viktor Bengs, Eyke Hüllermeier
- abstract@[open-review](#): The reliable identification of the “best” arm while keeping the sample complexity as low as possible is a common task in the field of multi-armed bandits. In the multi-dueling variant of multi-armed bandits, where feedback is provided in the form of a winning arm among a set of k chosen ones, a reasonable notion of best arm is the generalized Condorcet winner (GCW). The latter is the arm that has the greatest probability of being the winner in each subset containing it. In this paper, we derive lower bounds on the sample complexity for the task of identifying the GCW under various assumptions. As a by-product, our lower bound results provide new insights for the special case of dueling bandits ($k = 2$). We propose the Dvoretzky–Kiefer–Wolfowitz tournament (DKWT) algorithm, which we prove to be nearly optimal. In a numerical study, we show that DKWT empirically outperforms current state-of-the-art algorithms, even in the special case of dueling bandits or under a Plackett–Luce assumption on the feedback mechanism.

[Robust Inverse Reinforcement Learning under Transition Dynamics Mismatch](#)

- Luca Viano, Yu-Ting Huang, Parameswaran Kamalaruban, Adrian Weller, Volkan Cevher
- abstract@[open-review](#): We study the inverse reinforcement learning (IRL) problem under a transition dynamics mismatch between the expert and the learner. Specifically, we consider the Maximum Causal Entropy (MCE) IRL learner model and provide a tight upper bound on the learner's performance degradation based on the ℓ_1 -distance between the transition dynamics of the expert and the learner. Leveraging insights from the Robust RL literature, we propose a robust MCE IRL algorithm, which is a principled approach to help with this mismatch. Finally, we empirically demonstrate the stable performance of our algorithm compared to the standard MCE IRL algorithm under transition dynamics mismatches in both finite and continuous MDP problems.

[Re-ranking for image retrieval and transductive few-shot classification](#)

- Xi SHEN, Yang Xiao, Shell Xu Hu, Othman Sbai, Mathieu Aubry
- abstract@[open-review](#): In the problems of image retrieval and few-shot classification, the mainstream approaches focus on learning a better feature representation. However, directly tackling the distance or similarity measure between images could also be efficient. To this end, we revisit the idea of re-ranking the top-k retrieved images in the context of image retrieval (e.g., the k-reciprocal nearest neighbors) and generalize this idea to transductive few-shot learning. We propose to meta-learn the re-ranking updates such that the similarity graph converges towards the target similarity graph induced by the image labels. Specifically, the re-ranking module takes as input an initial similarity graph between the query image and the contextual images using a pre-trained feature extractor, and predicts an improved similarity graph by leveraging the structure among the involved images. We show that our re-ranking approach can be applied to unseen images and can further boost existing approaches for both image retrieval and few-shot learning problems. Our approach operates either independently or in conjunction with classical re-ranking approaches, yielding clear and consistent improvements on image retrieval (CUB, Cars, SOP, rOxford5K and rParis6K) and transductive few-shot classification (Mini-ImageNet, tiered-ImageNet and CIFAR-FS) benchmarks. Our code is available at <https://imagine.enpc.fr/~shenx/SSR/>.

[Post-processing for Individual Fairness](#)

- Felix Petersen, Debarghya Mukherjee, Yuekai Sun, Mikhail Yurochkin
- abstract@[open-review](#): Post-processing in algorithmic fairness is a versatile approach for correcting bias in ML systems that are already used in production. The main appeal of post-processing is that it avoids expensive retraining. In this work, we propose general post-processing algorithms for individual fairness (IF). We consider a setting where the learner only has access to the predictions of the original model and a similarity graph between individuals, guiding the desired fairness constraints. We cast the IF post-processing problem as a graph smoothing problem corresponding to graph Laplacian regularization that preserves the desired "treat similar individuals similarly" interpretation. Our theoretical results demonstrate the connection of the new objective function to a local relaxation of the original individual fairness. Empirically, our post-processing algorithms correct individual biases in large-scale NLP models such as BERT, while preserving accuracy.

OpenMatch: Open-Set Semi-supervised Learning with Open-set Consistency Regularization

- Kuniaki Saito, Donghyun Kim, Kate Saenko
- abstract@[open-review](#): Semi-supervised learning (SSL) is an effective means to leverage unlabeled data to improve a model's performance. Typical SSL methods like FixMatch assume that labeled and unlabeled data share the same label space. However, in practice, unlabeled data can contain categories unseen in the labeled set, i.e., outliers, which can significantly harm the performance of SSL algorithms. To address this problem, we propose a novel Open-set Semi-Supervised Learning (OSSL) approach called OpenMatch. Learning representations of inliers while rejecting outliers is essential for the success of OSSL. To this end, OpenMatch unifies FixMatch with novelty detection based on one-vs-all (OVA) classifiers. The OVA-classifier outputs the confidence score of a sample being an inlier, providing a threshold to detect outliers. Another key contribution is an open-set soft-consistency regularization loss, which enhances the smoothness of the OVA-classifier with respect to input transformations and greatly improves outlier detection. Ours achieves state-of-the-art performance on three datasets, and even outperforms a fully supervised model in detecting outliers unseen in unlabeled data on CIFAR10. The code is available at \url{https://github.com/VisionLearningGroup/OP_Match}.

End-to-End Training of Multi-Document Reader and Retriever for Open-Domain Question Answering

- Devendra Singh, Siva Reddy, Will Hamilton, Chris Dyer, Dani Yogatama
- abstract@[open-review](#): We present an end-to-end differentiable training method for retrieval-augmented open-domain question answering systems that combine information from multiple retrieved documents when generating answers. We model retrieval decisions as latent variables over sets of relevant documents. Since marginalizing over sets of retrieved documents is computationally hard, we approximate this using an expectation-maximization algorithm. We iteratively estimate the value of our latent variable (the set of relevant documents for a given question) and then use this estimate to update the retriever and reader parameters. We hypothesize that such end-to-end training allows training signals to flow to the reader and then to the retriever better than staged-wise training. This results in a retriever that is able to select more relevant documents for a question and a reader that is trained on more accurate documents to generate an answer. Experiments on three benchmark datasets demonstrate that our proposed method outperforms all existing approaches of comparable size by 2-3% absolute exact match points, achieving new state-of-the-art results. Our results also demonstrate the feasibility of learning to retrieve to improve answer generation without explicit supervision of retrieval decisions.

Fast Algorithms for L_{∞} -constrained S-rectangular Robust MDPs

- Bahram Behzadian, Marek Petrik, Chin Pang Ho
- abstract@[open-review](#): Robust Markov decision processes (RMDPs) are a useful building block of robust reinforcement learning algorithms but can be hard to solve. This paper proposes a fast, exact algorithm for computing the Bellman operator for S-rectangular robust Markov decision processes with L_{∞} -constrained rectangular ambiguity sets. The algorithm combines a novel homotopy continuation method with a bisection method to solve S-rectangular ambiguity in quasi-linear time in the number of states and actions. The algorithm improves on the cubic time required by leading general linear programming methods. Our experimental results confirm the practical viability of our method and show that it outperforms a leading commercial optimization package by several orders of magnitude.

Instance-optimal Mean Estimation Under Differential Privacy

- Ziyue Huang, Yuting Liang, Ke Yi
- abstract@[open-review](#): Mean estimation under differential privacy is a fundamental problem, but worst-case optimal mechanisms do not offer meaningful utility guarantees in practice when the global sensitivity is very large. Instead, various heuristics have been proposed to reduce the error on real-world data that do not resemble the worst-case instance. This paper takes a principled approach, yielding a mechanism that is instance-optimal in a strong sense. In addition to its theoretical optimality, the mechanism is also simple and practical, and adapts to a variety of data characteristics without the need of parameter tuning. It easily extends to the local and shuffle model as well.

Look at the Variance! Efficient Black-box Explanations with Sobol-based Sensitivity Analysis

- Thomas FEL, Remi Cadene, Mathieu Chalvidal, Matthieu Cord, David Vigouroux, Thomas Serre
- abstract@[open-review](#): We describe a novel attribution method which is grounded in Sensitivity Analysis and uses Sobol indices. Beyond modeling the individual contributions of image regions, Sobol indices provide an efficient way to capture higher-order interactions between image regions and their contributions to a neural network's prediction through the lens of variance. We describe an approach that makes the computation of these indices efficient for high-dimensional problems by using perturbation masks coupled with efficient estimators to handle the high dimensionality of images. Importantly, we show that the proposed method leads to favorable scores on standard benchmarks for vision (and language models) while drastically reducing the computing time compared to other black-box methods -- even surpassing the accuracy of state-of-the-art white-box methods which require access to internal representations. Our code is freely available: github.com/fel-thomas/Sobol-Attribution-Method.

PatchGame: Learning to Signal Mid-level Patches in Referential Games

- Kamal Gupta, Gowthami Somepalli, Anubhav Anubhav, Vinoj Yasanga Jayasundara Magalle Hewa, Matthias Zwicker, Abhinav Shrivastava
- abstract@[open-review](#): We study a referential game (a type of signaling game) where two agents communicate with each other via a discrete bottleneck to achieve a common goal. In our referential game, the goal of the speaker is to compose a message or a symbolic representation of "important" image patches, while the task for the listener is to match the speaker's message to a different view of the same image. We show that it is indeed possible for the two agents to develop a communication protocol without explicit or implicit supervision. We further investigate the developed protocol and show the applications in speeding up recent Vision Transformers by using only important patches, and as pre-training for downstream recognition tasks (e.g., classification).

Implicit Generative Copulas

- Tim Janke, Mohamed Ghanmi, Florian Steinke
- abstract@[open-review](#): Copulas are a powerful tool for modeling multivariate distributions as they allow to separately estimate the univariate marginal distributions and the joint dependency structure. However, known parametric copulas offer limited flexibility especially in high dimensions, while commonly used non-parametric methods suffer from the curse of dimensionality. A popular remedy is to construct a tree-based hierarchy of conditional

bivariate copulas. In this paper, we propose a flexible, yet conceptually simple alternative based on implicit generative neural networks. The key challenge is to ensure marginal uniformity of the estimated copula distribution. We achieve this by learning a multivariate latent distribution with unspecified marginals but the desired dependency structure. By applying the probability integral transform, we can then obtain samples from the high-dimensional copula distribution without relying on parametric assumptions or the need to find a suitable tree structure. Experiments on synthetic and real data from finance, physics, and image generation demonstrate the performance of this approach.

[Tensor Normal Training for Deep Learning Models](#)

- Yi Ren, Donald Goldfarb
- abstract@[open-review](#): Despite the predominant use of first-order methods for training deep learning models, second-order methods, and in particular, natural gradient methods, remain of interest because of their potential for accelerating training through the use of curvature information. Several methods with non-diagonal preconditioning matrices, including KFAC, Shampoo, and K-BFGS, have been proposed and shown to be effective. Based on the so-called tensor normal (TN) distribution, we propose and analyze a brand new approximate natural gradient method, Tensor Normal Training (TNT), which like Shampoo, only requires knowledge of the shape of the training parameters. By approximating the probabilistically based Fisher matrix, as opposed to the empirical Fisher matrix, our method uses the block-wise covariance of the sampling based gradient as the pre-conditioning matrix. Moreover, the assumption that the sampling-based (tensor) gradient follows a TN distribution, ensures that its covariance has a Kronecker separable structure, which leads to a tractable approximation to the Fisher matrix. Consequently, TNT's memory requirements and per-iteration computational costs are only slightly higher than those for first-order methods. In our experiments, TNT exhibited superior optimization performance to state-of-the-art first-order methods, and comparable optimization performance to the state-of-the-art second-order methods KFAC and Shampoo. Moreover, TNT demonstrated its ability to generalize as well as first-order methods, while using fewer epochs.

[Unintended Selection: Persistent Qualification Rate Disparities and Interventions](#)

- Reilly Raab, Yang Liu
- abstract@[open-review](#): Realistically---and equitably---modeling the dynamics of group-level disparities in machine learning remains an open problem. In particular, we desire models that do not suppose inherent differences between artificial groups of people---but rather endogenize disparities by appeal to unequal initial conditions of insular subpopulations. In this paper, agents each have a real-valued feature $\$X\$$ (e.g., credit score) informed by a ``true'' binary label $\$Y\$$ representing qualification (e.g., for a loan). Each agent alternately (1) receives a binary classification label $\$hat{Y}\$$ (e.g., loan approval) from a Bayes-optimal machine learning classifier observing $\$X\$$ and (2) may update their qualification $\$Y\$$ by imitating successful strategies (e.g., seek a raise) within an isolated group $\$G\$$ of agents to which they belong. We consider the disparity of qualification rates $\$Pr(Y=1)\$$ between different groups and how this disparity changes subject to a sequence of Bayes-optimal classifiers repeatedly retrained on the global population. We model the evolving qualification rates of each subpopulation (group) using the replicator equation, which derives from a class of imitation processes. We show that differences in qualification rates between subpopulations can persist indefinitely for a set of non-trivial equilibrium states due to uniformed classifier deployments, even when groups are identical in all aspects except initial qualification densities. We next simulate the effects of commonly proposed fairness interventions on this dynamical system along with a new feedback control mechanism capable of permanently eliminating group-level qualification rate disparities. We conclude by discussing the limitations of our model and findings and by outlining potential future work.

[Revisiting 3D Object Detection From an Egocentric Perspective](#)

- Boyang Deng, Charles R Qi, Mahyar Najibi, Thomas Funkhouser, Yin Zhou, Dragomir Anguelov
- abstract@[open-review](#): 3D object detection is a key module for safety-critical robotics applications such as autonomous driving. For these applications, we care most about how the detections affect the ego-agent's behavior and safety (the egocentric perspective). Intuitively, we seek more accurate descriptions of object geometry when it's more likely to interfere with the ego-agent's motion trajectory. However, current detection metrics, based on box Intersection-over-Union (IoU), are object-centric and aren't designed to capture the spatio-temporal relationship between objects and the ego-agent. To address this issue, we propose a new egocentric measure to evaluate 3D object detection, namely Support Distance Error (SDE). Our analysis based on SDE reveals that the egocentric detection quality is bounded by the coarse geometry of the bounding boxes. Given the insight that SDE would benefit from more accurate geometry descriptions, we propose to represent objects as amodal contours, specifically amodal star-shaped polygons, and devise a simple model, StarPoly, to predict such contours. Our experiments on the large-scale Waymo Open Dataset show that SDE better reflects the impact of detection quality on the ego-agent's safety compared to IoU; and the estimated contours from StarPoly consistently improve the egocentric detection quality over recent 3D object detectors.

[Optimizing Information-theoretical Generalization Bound via Anisotropic Noise of SGLD](#)

- Bohan Wang, Huishuai Zhang, Jieyu Zhang, Qi Meng, Wei Chen, Tie-Yan Liu
- abstract@[open-review](#): Recently, the information-theoretical framework has been proven to be able to obtain non-vacuous generalization bounds for large models trained by Stochastic Gradient Langevin Dynamics (SGLD) with isotropic noise. In this paper, we optimize the information-theoretical generalization bound by manipulating the noise structure in SGLD. We prove that with constraint to guarantee low empirical risk, the optimal noise covariance is the square root of the expected gradient covariance if both the prior and the posterior are jointly optimized. This validates that the optimal noise is quite close to the empirical gradient covariance. Technically, we develop a new information-theoretical bound that enables such an optimization analysis. We then apply matrix analysis to derive the form of optimal noise covariance. Presented constraint and results are validated by the empirical observations.

[Addressing Algorithmic Disparity and Performance Inconsistency in Federated Learning](#)

- Sen Cui, Weishen Pan, Jian Liang, Changshui Zhang, Fei Wang
- abstract@[open-review](#): Federated learning (FL) has gained growing interests for its capability of learning from distributed data sources collectively without the need of accessing the raw data samples across different sources. So far FL research has mostly focused on improving the performance, how the algorithmic disparity will be impacted for the model learned from FL and the impact of algorithmic disparity on the utility inconsistency are largely unexplored. In this paper, we propose an FL framework to jointly consider performance consistency and algorithmic fairness across different local clients (data sources). We derive our framework from a constrained multi-objective optimization perspective, in which we learn a model satisfying fairness constraints on all clients with consistent performance. Specifically, we treat the algorithm prediction loss at each local client as an objective and maximize the worst-performing client with fairness constraints through optimizing a surrogate maximum function with all objectives involved. A gradient-based procedure is employed to achieve the Pareto optimality of this optimization problem. Theoretical analysis is provided to prove that our method can converge to a Pareto solution that achieves the min-max performance with fairness constraints on all clients. Comprehensive experiments on synthetic and real-world datasets demonstrate the superiority of our approach over baselines and its effectiveness in achieving both fairness and consistency across all local clients.

[A Mathematical Framework for Quantifying Transferability in Multi-source Transfer Learning](#)

- Xinyi Tong, Xiangxiang Xu, Shao-Lun Huang, Lizhong Zheng
- abstract@[open-review](#): Current transfer learning algorithm designs mainly focus on the similarities between source and target tasks, while the impacts of the sample sizes of these tasks are often not sufficiently addressed. This paper proposes a mathematical framework for quantifying the transferability in

multi-source transfer learning problems, with both the task similarities and the sample complexity of learning models taken into account. In particular, we consider the setup where the models learned from different tasks are linearly combined for learning the target task, and use the optimal combining coefficients to measure the transferability. Then, we demonstrate the analytical expression of this transferability measure, characterized by the sample sizes, model complexity, and the similarities between source and target tasks, which provides fundamental insights of the knowledge transferring mechanism and the guidance for algorithm designs. Furthermore, we apply our analyses for practical learning tasks, and establish a quantifiable transferability measure by exploiting a parameterized model. In addition, we develop an alternating iterative algorithm to implement our theoretical results for training deep neural networks in multi-source transfer learning tasks. Finally, experiments on image classification tasks show that our approach outperforms existing transfer learning algorithms in multi-source and few-shot scenarios.

[Morié Attack \(MA\): A New Potential Risk of Screen Photos](#)

- Dantong Niu, Ruohao Guo, Yisen Wang
- abstract@[open-review](#): Images, captured by a camera, play a critical role in training Deep Neural Networks (DNNs). Usually, we assume the images acquired by cameras are consistent with the ones perceived by human eyes. However, due to the different physical mechanisms between human-vision and computer-vision systems, the final perceived images could be very different in some cases, for example shooting on digital monitors. In this paper, we find a special phenomenon in digital image processing, the moiré effect, that could cause unnoticed security threats to DNNs. Based on it, we propose a Moiré Attack (MA) that generates the physical-world moiré pattern adding to the images by mimicking the shooting process of digital devices. Extensive experiments demonstrate that our proposed digital Moiré Attack (MA) is a perfect camouflage for attackers to tamper with DNNs with a high success rate (\$100.0\%\$ for untargeted and \$97.0\%\$ for targeted attack with the noise budget \$\epsilon=4\$), high transferability rate across different models, and high robustness under various defenses. Furthermore, MA owns great stealthiness because the moiré effect is unavoidable due to the camera's inner physical structure, which therefore hardly attracts the awareness of humans. Our code is available at https://github.com/Dantong88/Moire_Attack.

[Fast Bayesian Inference for Gaussian Cox Processes via Path Integral Formulation](#)

- Hideaki Kim
- abstract@[open-review](#): Gaussian Cox processes are widely-used point process models that use a Gaussian process to describe the Bayesian a priori uncertainty present in latent intensity functions. In this paper, we propose a novel Bayesian inference scheme for Gaussian Cox processes by exploiting a conceptually-intuitive \$\{Y_t\}\$ path integral formulation. The proposed scheme does not rely on domain discretization, scales linearly with the number of observed events, has a lower complexity than the state-of-the-art variational Bayesian schemes with respect to the number of inducing points, and is applicable to a wide range of Gaussian Cox processes with various types of link functions. Our scheme is especially beneficial under the multi-dimensional input setting, where the number of inducing points tends to be large. We evaluate our scheme on synthetic and real-world data, and show that it achieves comparable predictive accuracy while being tens of times faster than reference methods.

[Lattice partition recovery with dyadic CART](#)

- OSCAR HERNAN MADRID PADILLA, Yi Yu, Alessandro Rinaldo
- abstract@[open-review](#): We study piece-wise constant signals corrupted by additive Gaussian noise over a d -dimensional lattice. Data of this form naturally arise in a host of applications, and the tasks of signal detection or testing, de-noising and estimation have been studied extensively in the statistical and signal processing literature. In this paper we consider instead the problem of partition recovery, i.e., of estimating the partition of the lattice induced by the constancy regions of the unknown signal, using the computationally-efficient dyadic classification and regression tree (DCART) methodology proposed by \citet{donoho1997cart}. We prove that, under appropriate regularity conditions on the shape of the partition elements, a DCART-based procedure consistently estimates the underlying partition at a rate of order $\sigma^2 k \log(N)/\kappa^2$, where k is the minimal number of rectangular sub-graphs obtained using recursive dyadic partitions supporting the signal partition, σ^2 is the noise variance, κ is the minimal magnitude of the signal difference among contiguous elements of the partition and N is the size of the lattice. Furthermore, under stronger assumptions, our method attains a sharper estimation error of order $\sigma^2 \log(N)/\kappa^2$, independent of k , which we show to be minimax rate optimal. Our theoretical guarantees further extend to the partition estimator based on the optimal regression tree estimator (ORT) of \citet{chatterjee2019adaptive} and to the one obtained through an NP-hard exhaustive search method. We corroborate our theoretical findings and the effectiveness of DCART for partition recovery in simulations.

[Robust Deep Reinforcement Learning through Adversarial Loss](#)

- Tuomas Oikarinen, Wang Zhang, Alexandre Megretski, Luca Daniel, Tsui-Wei Weng
- abstract@[open-review](#): Recent studies have shown that deep reinforcement learning agents are vulnerable to small adversarial perturbations on the agent's inputs, which raises concerns about deploying such agents in the real world. To address this issue, we propose RADIAL-RL, a principled framework to train reinforcement learning agents with improved robustness against L_p -norm bounded adversarial attacks. Our framework is compatible with popular deep reinforcement learning algorithms and we demonstrate its performance with deep Q-learning, A3C and PPO. We experiment on three deep RL benchmarks (Atari, MuJoCo and ProcGen) to show the effectiveness of our robust training algorithm. Our RADIAL-RL agents consistently outperform prior methods when tested against attacks of varying strength and are more computationally efficient to train. In addition, we propose a new evaluation method called Greedy-Worst-Case Reward (GWC) to measure attack agnostic robustness of deep RL agents. We show that GWC can be evaluated efficiently and is a good estimate of the reward under the worst possible sequence of adversarial attacks. All code used for our experiments is available at https://github.com/tuomaso/radial_rl_v2.

[Provable Model-based Nonlinear Bandit and Reinforcement Learning: Shelve Optimism, Embrace Virtual Curvature](#)

- Kefan Dong, Jiaqi Yang, Tengyu Ma
- abstract@[open-review](#): This paper studies model-based bandit and reinforcement learning (RL) with nonlinear function approximations. We propose to study convergence to approximate local maxima because we show that global convergence is statistically intractable even for one-layer neural net bandit with a deterministic reward. For both nonlinear bandit and RL, the paper presents a model-based algorithm, Virtual Ascent with Online Model Learner (ViOlin), which provably converges to a local maximum with sample complexity that only depends on the sequential Rademacher complexity of the model class. Our results imply novel global or local regret bounds on several concrete settings such as linear bandit with finite or sparse model class, and two-layer neural net bandit. A key algorithmic insight is that optimism may lead to over-exploration even for two-layer neural net model class. On the other hand, for convergence to local maxima, it suffices to maximize the virtual return if the model can also reasonably predict the gradient and Hessian of the real return.

[You Only Look at One Sequence: Rethinking Transformer in Vision through Object Detection](#)

- Yuxin Fang, Bencheng Liao, Xinggang Wang, Jiemin Fang, Jiyang Qi, Rui Wu, Jianwei Niu, Wenyu Liu
- abstract@[open-review](#): Can Transformer perform D -object- and region-level recognition from a pure sequence-to-sequence perspective with minimal knowledge about the D spatial structure? To answer this question, we present You Only Look at One Sequence (YOLOS), a series of object detection models based on the vanilla Vision Transformer with the fewest possible modifications, region priors, as well as inductive biases of the target task. We find that YOLOS pre-trained on the mid-sized ImageNet-1k dataset only can already achieve quite competitive performance on the challenging COCO object detection benchmark, e.g., YOLOS-Base directly adopted from BERT-Base architecture can obtain 42.0 box AP on

COCO val. We also discuss the impacts as well as limitations of current pre-train schemes and model scaling strategies for Transformer in vision through YOLOS. Code and pre-trained models are available at <https://github.com/hustvl/YOLOS>.

[Learning to delegate for large-scale vehicle routing](#)

- Sirui Li, Zhongxia Yan, Cathy Wu
- abstract@[open-review](#): Vehicle routing problems (VRPs) form a class of combinatorial problems with wide practical applications. While previous heuristic or learning-based works achieve decent solutions on small problem instances, their performance deteriorates in large problems. This article presents a novel learning-augmented local search framework to solve large-scale VRP. The method iteratively improves the solution by identifying appropriate subproblems and \$delegating\$ their improvement to a black box subsolver. At each step, we leverage spatial locality to consider only a linear number of subproblems, rather than exponential. We frame subproblem selection as regression and train a Transformer on a generated training set of problem instances. Our method accelerates state-of-the-art VRP solvers by 10x to 100x while achieving competitive solution qualities for VRPs with sizes ranging from 500 to 3000. Learned subproblem selection offers a 1.5x to 2x speedup over heuristic or random selection. Our results generalize to a variety of VRP distributions, variants, and solvers.

[Effective Meta-Regularization by Kernelized Proximal Regularization](#)

- Weisen Jiang, James Kwok, Yu Zhang
- abstract@[open-review](#): We study the problem of meta-learning, which has proved to be advantageous to accelerate learning new tasks with a few samples. The recent approaches based on deep kernels achieve the state-of-the-art performance. However, the regularizers in their base learners are not learnable. In this paper, we propose an algorithm called MetaProx to learn a proximal regularizer for the base learner. We theoretically establish the convergence of MetaProx. Experimental results confirm the advantage of the proposed algorithm.

[Towards Context-Agnostic Learning Using Synthetic Data](#)

- Charles Jin, Martin Rinard
- abstract@[open-review](#): We propose a novel setting for learning, where the input domain is the image of a map defined on the product of two sets, one of which completely determines the labels. We derive a new risk bound for this setting that decomposes into a bias and an error term, and exhibits a surprisingly weak dependence on the true labels. Inspired by these results, we present an algorithm aimed at minimizing the bias term by exploiting the ability to sample from each set independently. We apply our setting to visual classification tasks, where our approach enables us to train classifiers on datasets that consist entirely of a single synthetic example of each class. On several standard benchmarks for real-world image classification, we achieve robust performance in the context-agnostic setting, with good generalization to real world domains, whereas training directly on real world data without our techniques yields classifiers that are brittle to perturbations of the background.

[Minimax Optimal Quantile and Semi-Adversarial Regret via Root-Logarithmic Regularizers](#)

- Jeffrey Negrea, Blair Bilodeau, Nicolò Campolongo, Francesco Orabona, Dan Roy
- abstract@[open-review](#): Quantile (and, more generally, KL) regret bounds, such as those achieved by NormalHedge (Chaudhuri, Freund, and Hsu 2009) and its variants, relax the goal of competing against the best individual expert to only competing against a majority of experts on adversarial data. More recently, the semi-adversarial paradigm (Bilodeau, Negrea, and Roy 2020) provides an alternative relaxation of adversarial online learning by considering data that may be neither fully adversarial nor stochastic (I.I.D.). We achieve the minimax optimal regret in both paradigms using FTRL with separate, novel, root-logarithmic regularizers, both of which can be interpreted as yielding variants of NormalHedge. We extend existing KL regret upper bounds, which hold uniformly over target distributions, to possibly uncountable expert classes with arbitrary priors; provide the first full-information lower bounds for quantile regret on finite expert classes (which are tight); and provide an adaptively minimax optimal algorithm for the semi-adversarial paradigm that adapts to the true, unknown constraint faster, leading to uniformly improved regret bounds over existing methods.

[Gradient-Free Adversarial Training Against Image Corruption for Learning-based Steering](#)

- Yu Shen, Laura Zheng, Manli Shu, Weizi Li, Tom Goldstein, Ming Lin
- abstract@[open-review](#): We introduce a simple yet effective framework for improving the robustness of learning algorithms against image corruptions for autonomous driving. These corruptions can occur due to both internal (e.g., sensor noises and hardware abnormalities) and external factors (e.g., lighting, weather, visibility, and other environmental effects). Using sensitivity analysis with FID-based parameterization, we propose a novel algorithm exploiting basis perturbations to improve the overall performance of autonomous steering and other image processing tasks, such as classification and detection, for self-driving cars. Our model not only improves the performance on the original dataset, but also achieves significant performance improvement on datasets with multiple and unseen perturbations, up to 87% and 77%, respectively. A comparison between our approach and other SOTA techniques confirms the effectiveness of our technique in improving the robustness of neural network training for learning-based steering and other image processing tasks.

[Deep Proxy Causal Learning and its Application to Confounded Bandit Policy Evaluation](#)

- Liyuan Xu, Heishiro Kanagawa, Arthur Gretton
- abstract@[open-review](#): Proxy causal learning (PCL) is a method for estimating the causal effect of treatments on outcomes in the presence of unobserved confounding, using proxies (structured side information) for the confounder. This is achieved via two-stage regression: in the first stage, we model relations among the treatment and proxies; in the second stage, we use this model to learn the effect of treatment on the outcome, given the context provided by the proxies. PCL guarantees recovery of the true causal effect, subject to identifiability conditions. We propose a novel method for PCL, the deep feature proxy variable method (DFPV), to address the case where the proxies, treatments, and outcomes are high-dimensional and have nonlinear complex relationships, as represented by deep neural network features. We show that DFPV outperforms recent state-of-the-art PCL methods on challenging synthetic benchmarks, including settings involving high dimensional image data. Furthermore, we show that PCL can be applied to off-policy evaluation for the confounded bandit problem, in which DFPV also exhibits competitive performance.

[Certifying Robustness to Programmable Data Bias in Decision Trees](#)

- Anna Meyer, Aws Albargouthi, Loris D'Antoni
- abstract@[open-review](#): Datasets can be biased due to societal inequities, human biases, under-representation of minorities, etc. Our goal is to certify that models produced by a learning algorithm are pointwise-robust to dataset biases. This is a challenging problem: it entails learning models for a large, or even infinite, number of datasets, ensuring that they all produce the same prediction. We focus on decision-tree learning due to the interpretable nature of the models. Our approach allows programmatically specifying \emph{bias models} across a variety of dimensions (e.g., label-flipping or missing data), composing types of bias, and targeting bias towards a specific group. To certify robustness, we use a novel symbolic technique to evaluate a decision-tree learner on a large, or infinite, number of datasets, certifying that each and every dataset produces the same prediction for a specific test point. We evaluate our approach on datasets that are commonly used in the fairness literature, and demonstrate our approach's viability on a range of bias models.

[TöRF: Time-of-Flight Radiance Fields for Dynamic Scene View Synthesis](#)

- Benjamin Attal, Eliot Laidlaw, Aaron Gokaslan, Changil Kim, Christian Richardt, James Tompkin, Matthew O'Toole
- abstract@[open-review](#): Neural networks can represent and accurately reconstruct radiance fields for static 3D scenes (e.g., NeRF). Several works extend these to dynamic scenes captured with monocular video, with promising performance. However, the monocular setting is known to be an under-constrained problem, and so methods rely on data-driven priors for reconstructing dynamic content. We replace these priors with measurements from a time-of-flight (ToF) camera, and introduce a neural representation based on an image formation model for continuous-wave ToF cameras. Instead of working with processed depth maps, we model the raw ToF sensor measurements to improve reconstruction quality and avoid issues with low reflectance regions, multi-path interference, and a sensor's limited unambiguous depth range. We show that this approach improves robustness of dynamic scene reconstruction to erroneous calibration and large motions, and discuss the benefits and limitations of integrating RGB+ToF sensors now available on modern smartphones.

[Sequence-to-Sequence Learning with Latent Neural Grammars](#)

- Yoon Kim
- abstract@[open-review](#): Sequence-to-sequence learning with neural networks has become the de facto standard for sequence modeling. This approach typically models the local distribution over the next element with a powerful neural network that can condition on arbitrary context. While flexible and performant, these models often require large datasets for training and can fail spectacularly on benchmarks designed to test for compositional generalization. This work explores an alternative, hierarchical approach to sequence-to-sequence learning with synchronous grammars, where each node in the target tree is transduced by a subset of nodes in the source tree. The source and target trees are treated as fully latent and marginalized out during training. We develop a neural parameterization of the grammar which enables parameter sharing over combinatorial structures without the need for manual feature engineering. We apply this latent neural grammar to various domains---a diagnostic language navigation task designed to test for compositional generalization (SCAN), style transfer, and small-scale machine translation---and find that it performs respectably compared to standard baselines.

[Exploration-Exploitation in Multi-Agent Competition: Convergence with Bounded Rationality](#)

- Stefanos Leonardos, Georgios Piliouras, Kelly Spendlove
- abstract@[open-review](#): The interplay between exploration and exploitation in competitive multi-agent learning is still far from being well understood. Motivated by this, we study smooth Q-learning, a prototypical learning model that explicitly captures the balance between game rewards and exploration costs. We show that Q-learning always converges to the unique quantal-response equilibrium (QRE), the standard solution concept for games under bounded rationality, in weighted zero-sum polymatrix games with heterogeneous learning agents using positive exploration rates. Complementing recent results about convergence in weighted potential games [16,34], we show that fast convergence of Q-learning in competitive settings obtains regardless of the number of agents and without any need for parameter fine-tuning. As showcased by our experiments in network zero-sum games, these theoretical results provide the necessary guarantees for an algorithmic approach to the currently open problem of equilibrium selection in competitive multi-agent settings.

[Low-Rank Extragradient Method for Nonsmooth and Low-Rank Matrix Optimization Problems](#)

- Atara Kaplan, Dan Garber
- abstract@[open-review](#): Low-rank and nonsmooth matrix optimization problems capture many fundamental tasks in statistics and machine learning. While significant progress has been made in recent years in developing efficient methods for \textit{smooth} low-rank optimization problems that avoid maintaining high-rank matrices and computing expensive high-rank SVDs, advances for nonsmooth problems have been slow paced. In this paper we consider standard convex relaxations for such problems. Mainly, we prove that under a natural \textit{generalized strict complementarity} condition and under the relatively mild assumption that the nonsmooth objective can be written as a maximum of smooth functions, the \textit{extragradient method}, when initialized with a "warm-start" point, converges to an optimal solution with rate $\$O(1/t)$ while requiring only two \textit{low-rank} SVDs per iteration. We give a precise trade-off between the rank of the SVDs required and the radius of the ball in which we need to initialize the method. We support our theoretical results with empirical experiments on several nonsmooth low-rank matrix recovery tasks, demonstrating that using simple initializations, the extragradient method produces exactly the same iterates when full-rank SVDs are replaced with SVDs of rank that matches the rank of the (low-rank) ground-truth matrix to be recovered.

[Which Mutual-Information Representation Learning Objectives are Sufficient for Control?](#)

- Kate Rakelly, Abhishek Gupta, Carlos Florensa, Sergey Levine
- abstract@[open-review](#): Mutual information (MI) maximization provides an appealing formalism for learning representations of data. In the context of reinforcement learning (RL), such representations can accelerate learning by discarding irrelevant and redundant information, while retaining the information necessary for control. Much prior work on these methods has addressed the practical difficulties of estimating MI from samples of high-dimensional observations, while comparatively less is understood about which MI objectives yield representations that are sufficient for RL from a theoretical perspective. In this paper, we formalize the sufficiency of a state representation for learning and representing the optimal policy, and study several popular MI based objectives through this lens. Surprisingly, we find that two of these objectives can yield insufficient representations given mild and common assumptions on the structure of the MDP. We corroborate our theoretical results with empirical experiments on a simulated game environment with visual observations.

[A Geometric Perspective towards Neural Calibration via Sensitivity Decomposition](#)

- Junjiao Tian, Dylan Yung, Yen-Chang Hsu, Zsolt Kira
- abstract@[open-review](#): It is well known that vision classification models suffer from poor calibration in the face of data distribution shifts. In this paper, we take a geometric approach to this problem. We propose Geometric Sensitivity Decomposition (GSD) which decomposes the norm of a sample feature embedding and the angular similarity to a target classifier into an instance-dependent and an instance-independent component. The instance-dependent component captures the sensitive information about changes in the input while the instance-independent component represents the insensitive information serving solely to minimize the loss on the training dataset. Inspired by the decomposition, we analytically derive a simple extension to current softmax-linear models, which learns to disentangle the two components during training. On several common vision models, the disentangled model out-performs other calibration methods on standard calibration metrics in the face of out-of-distribution (OOD) data and corruption with significantly less complexity. Specifically, we surpass the current state of the art by 30.8% relative improvement on corrupted CIFAR100 in Expected Calibration Error.

[Towards a Unified Information-Theoretic Framework for Generalization](#)

- Mahdi Haghifam, Gintare Karolina Dziugaite, Shay Moran, Dan Roy
- abstract@[open-review](#): In this work, we investigate the expressiveness of the "conditional mutual information" (CMI) framework of Steinke and Zakynthinou (2020) and the prospect of using it to provide a unified framework for proving generalization bounds in the realizable setting. We first demonstrate that one can use this framework to express non-trivial (but sub-optimal) bounds for any learning algorithm that outputs hypotheses from a

class of bounded VC dimension. We then explore two directions of strengthening this bound: (i) Can the CMI framework express optimal bounds for VC classes? (ii) Can the CMI framework be used to analyze algorithms whose output hypothesis space is unrestricted (i.e. has an unbounded VC dimension)? With respect to Item (i) we prove that the CMI framework yields the optimal bound on the expected risk of Support Vector Machines (SVMs) for learning halfspaces. This result is an application of our general result showing that stable compression schemes Bousquet al. (2020) of size $\$k\$$ have uniformly bounded CMI of order $\$O(k)\$$. We further show that an inherent limitation of proper learning of VC classes contradicts the existence of a proper learner with constant CMI, and it implies a negative resolution to an open problem of Steinke and Zakynthinou (2020). We further study the CMI of empirical risk minimizers (ERMs) of class $\$H\$$ and show that it is possible to output all consistent classifiers (version space) with bounded CMI if and only if $\$H\$$ has a bounded star number (Hanneke and Yang (2015)). With respect to Item (ii) we prove a general reduction showing that "leave-one-out" analysis is expressible via the CMI framework. As a corollary we investigate the CMI of the one-inclusion-graph algorithm proposed by Haussler et al. (1994). More generally, we show that the CMI framework is universal in the sense that for every consistent algorithm and data distribution, the expected risk vanishes as the number of samples diverges if and only if its evaluated CMI has sublinear growth with the number of samples.

[Bayesian decision-making under misspecified priors with applications to meta-learning](#)

- Max Simchowitz, Christopher Tosh, Akshay Krishnamurthy, Daniel J. Hsu, Thodoris Lykouris, Miro Dudik, Robert E. Schapire
- abstract@[open-review](#): Thompson sampling and other Bayesian sequential decision-making algorithms are among the most popular approaches to tackle explore/exploit trade-offs in (contextual) bandits. The choice of prior in these algorithms offers flexibility to encode domain knowledge but can also lead to poor performance when misspecified. In this paper, we demonstrate that performance degrades gracefully with misspecification. We prove that the expected reward accrued by Thompson sampling (TS) with a misspecified prior differs by at most $\$tilde{O}(H^2 \epsilon)$ from TS with a well-specified prior, where ϵ is the total-variation distance between priors and H is the learning horizon. Our bound does not require the prior to have any parametric form. For priors with bounded support, our bound is independent of the cardinality or structure of the action space, and we show that it is tight up to universal constants in the worst case. Building on our sensitivity analysis, we establish generic PAC guarantees for algorithms in the recently studied Bayesian meta-learning setting and derive corollaries for various families of priors. Our results generalize along two axes: (1) they apply to a broader family of Bayesian decision-making algorithms, including a Monte-Carlo implementation of the knowledge gradient algorithm (KG), and (2) they apply to Bayesian POMDPs, the most general Bayesian decision-making setting, encompassing contextual bandits as a special case. Through numerical simulations, we illustrate how prior misspecification and the deployment of one-step look-ahead (as in KG) can impact the convergence of meta-learning in multi-armed and contextual bandits with structured and correlated priors.

[Neural Trees for Learning on Graphs](#)

- Rajat Talak, Siyi Hu, Lisa Peng, Luca Carlone
- abstract@[open-review](#): Graph Neural Networks (GNNs) have emerged as a flexible and powerful approach for learning over graphs. Despite this success, existing GNNs are constrained by their local message-passing architecture and are provably limited in their expressive power. In this work, we propose a new GNN architecture – the Neural Tree. The neural tree architecture does not perform message passing on the input graph, but on a tree-structured graph, called the H-tree, that is constructed from the input graph. Nodes in the H-tree correspond to subgraphs in the input graph, and they are reorganized in a hierarchical manner such that the parent of a node in the H-tree always corresponds to a larger subgraph in the input graph. We show that the neural tree architecture can approximate any smooth probability distribution function over an undirected graph. We also prove that the number of parameters needed to achieve an ϵ -approximation of the distribution function is exponential in the treewidth of the input graph, but linear in its size. We prove that any continuous G-invariant/equivariant function can be approximated by a nonlinear combination of such probability distribution functions over G. We apply the neural tree to semi-supervised node classification in 3D scene graphs, and show that these theoretical properties translate into significant gains in prediction accuracy, over the more traditional GNN architectures. We also show the applicability of the neural tree architecture to citation networks with large treewidth, by using a graph sub-sampling technique.

[Enabling Fast Differentially Private SGD via Just-in-Time Compilation and Vectorization](#)

- Pranav Subramani, Nicholas Vadivelu, Gautam Kamath
- abstract@[open-review](#): A common pain point in differentially private machine learning is the significant runtime overhead incurred when executing Differentially Private Stochastic Gradient Descent (DPSGD), which may be as large as two orders of magnitude. We thoroughly demonstrate that by exploiting powerful language primitives, including vectorization, just-in-time compilation, and static graph optimization, one can dramatically reduce these overheads, in many cases nearly matching the best non-private running times. These gains are realized in two frameworks: one is JAX, which provides rich support for these primitives through the XLA compiler. We also rebuild core parts of TensorFlow Privacy, integrating more effective vectorization as well as XLA compilation, granting significant memory and runtime improvements over previous release versions. Our proposed approaches allow us to achieve up to 50x speedups compared to the best alternatives. Our code is available at <https://github.com/TheSalon/fast-dpsgd>.

[The effectiveness of feature attribution methods and its correlation with automatic evaluation scores](#)

- Giang Nguyen, Daeyoung Kim, Anh Nguyen
- abstract@[open-review](#): Explaining the decisions of an Artificial Intelligence (AI) model is increasingly critical in many real-world, high-stake applications. Hundreds of papers have either proposed new feature attribution methods, discussed or harnessed these tools in their work. However, despite humans being the target end-users, most attribution methods were only evaluated on proxy automatic-evaluation metrics (Zhang et al. 2018; Zhou et al. 2016; Petsiuk et al. 2018). In this paper, we conduct the first user study to measure attribution map effectiveness in assisting humans in ImageNet classification and Stanford Dogs fine-grained classification, and when an image is natural or adversarial (i.e., contains adversarial perturbations). Overall, feature attribution is surprisingly not more effective than showing humans nearest training-set examples. On a harder task of fine-grained dog categorization, presenting attribution maps to humans does not help, but instead hurts the performance of human-AI teams compared to AI alone. Importantly, we found automatic attribution-map evaluation measures to correlate poorly with the actual human-AI team performance. Our findings encourage the community to rigorously test their methods on the downstream human-in-the-loop applications and to rethink the existing evaluation metrics.

[Coordinated Proximal Policy Optimization](#)

- Zifan Wu, Chao Yu, Deheng Ye, Junge Zhang, haiyin piao, Hankz Hankui Zhuo
- abstract@[open-review](#): We present Coordinated Proximal Policy Optimization (CoPPO), an algorithm that extends the original Proximal Policy Optimization (PPO) to the multi-agent setting. The key idea lies in the coordinated adaptation of step size during the policy update process among multiple agents. We prove the monotonicity of policy improvement when optimizing a theoretically-grounded joint objective, and derive a simplified optimization objective based on a set of approximations. We then interpret that such an objective in CoPPO can achieve dynamic credit assignment among agents, thereby alleviating the high variance issue during the concurrent update of agent policies. Finally, we demonstrate that CoPPO outperforms several strong baselines and is competitive with the latest multi-agent PPO method (i.e. MAPPO) under typical multi-agent settings, including cooperative matrix games and the StarCraft II micromanagement tasks.

[Unbiased Classification through Bias-Contrastive and Bias-Balanced Learning](#)

- Youngkyu Hong, Eunho Yang

- abstract@[open-review](#): Datasets for training machine learning models tend to be biased unless the data is collected with complete care. In such a biased dataset, models are susceptible to making predictions based on the biased features of the data. The biased model fails to generalize to the case where correlations between biases and targets are shifted. To mitigate this, we propose Bias-Contrastive (BiasCon) loss based on the contrastive learning framework, which effectively leverages the knowledge of bias labels. We further suggest Bias-Balanced (BiasBal) regression which trains the classification model toward the data distribution with balanced target-bias correlation. Furthermore, we propose Soft Bias-Contrastive (SoftCon) loss which handles the dataset without bias labels by softening the pair assignment of the BiasCon loss based on the distance in the feature space of the bias-capturing model. Our experiments show that our proposed methods significantly improve previous debiasing methods in various realistic datasets.

[Learning from Inside: Self-driven Siamese Sampling and Reasoning for Video Question Answering](#)

- Weijiang Yu, Haoteng Zheng, Mengfei Li, Lei Ji, Lijun Wu, Nong Xiao, Nan Duan
- abstract@[open-review](#): Recent advances in the video question answering (i.e., VideoQA) task have achieved strong success by following the paradigm of fine-tuning each clip-text pair independently on the pretrained transformer-based model via supervised learning. Intuitively, multiple samples (i.e., clips) should be interdependent to capture similar visual and key semantic information in the same video. To consider the interdependent knowledge between contextual clips into the network inference, we propose a Siamese Sampling and Reasoning (SiaSamRea) approach, which consists of a siamese sampling mechanism to generate sparse and similar clips (i.e., siamese clips) from the same video, and a novel reasoning strategy for integrating the interdependent knowledge between contextual clips into the network. The reasoning strategy contains two modules: (1) siamese knowledge generation to learn the inter-relationship among clips; (2) siamese knowledge reasoning to produce the refined soft label by propagating the weights of inter-relationship to the predicted candidates of all clips. Finally, our SiaSamRea can endow the current multimodal reasoning paradigm with the ability of learning from inside via the guidance of soft labels. Extensive experiments demonstrate our SiaSamRea achieves state-of-the-art performance on five VideoQA benchmarks, e.g., a significant +2.1% gain on MSRVTT-QA, +2.9% on MSVD-QA, +1.0% on ActivityNet-QA, +1.8% on How2QA and +4.3% (action) on TGIF-QA.

[Identification and Estimation of Joint Probabilities of Potential Outcomes in Observational Studies with Covariate Information](#)

- Ryusei Shingaki, manabu kuroki
- abstract@[open-review](#): The joint probabilities of potential outcomes are fundamental components of causal inference in the sense that (i) if they are identifiable, then the causal risk is also identifiable, but not vice versa (Pearl, 2009; Tian and Pearl, 2000) and (ii) they enable us to evaluate the probabilistic aspects of necessity", sufficiency", and ``necessity and sufficiency", which are important concepts of successful explanation (Watson, et al., 2020). However, because they are not identifiable without any assumptions, various assumptions have been utilized to evaluate the joint probabilities of potential outcomes, e.g., the assumption of monotonicity (Pearl, 2009; Tian and Pearl, 2000), the independence between potential outcomes (Robins and Richardson, 2011), the condition of gain equality (Li and Pearl, 2019), and the specific functional relationships between cause and effect (Pearl, 2009). Unlike existing identification conditions, in order to evaluate the joint probabilities of potential outcomes without such assumptions, this paper proposes two types of novel identification conditions using covariate information. In addition, when the joint probabilities of potential outcomes are identifiable through the proposed conditions, the estimation problem of the joint probabilities of potential outcomes reduces to that of singular models and thus they can not be evaluated by standard statistical estimation methods. To solve the problem, this paper proposes a new statistical estimation method based on the augmented Lagrangian method and shows the asymptotic normality of the proposed estimators. Given space constraints, the proofs, the details on the statistical estimation method, some numerical experiments, and the case study are provided in the supplementary material.

[Online false discovery rate control for anomaly detection in time series](#)

- Quentin Rebjock, Baris Kurt, Tim Januschowski, Laurent Callot
- abstract@[open-review](#): This article proposes novel rules for false discovery rate control (FDRC) geared towards online anomaly detection in time series. Online FDRC rules allow to control the properties of a sequence of statistical tests. In the context of anomaly detection, the null hypothesis is that an observation is normal and the alternative is that it is anomalous. FDRC rules allow users to target a lower bound on precision in unsupervised settings. The methods proposed in this article overcome short-comings of previous FDRC rules in the context of anomaly detection, in particular ensuring that power remains high even when the alternative is exceedingly rare (typical in anomaly detection) and the test statistics are serially dependent (typical in time series). We show the soundness of these rules in both theory and experiments.

[Pragmatic Image Compression for Human-in-the-Loop Decision-Making](#)

- Sid Reddy, Anca Dragan, Sergey Levine
- abstract@[open-review](#): Standard lossy image compression algorithms aim to preserve an image's appearance, while minimizing the number of bits needed to transmit it. However, the amount of information actually needed by the user for downstream tasks -- e.g., deciding which product to click on in a shopping website -- is likely much lower. To achieve this lower bitrate, we would ideally only transmit the visual features that drive user behavior, while discarding details irrelevant to the user's decisions. We approach this problem by training a compression model through human-in-the-loop learning as the user performs tasks with the compressed images. The key insight is to train the model to produce a compressed image that induces the user to take the same action that they would have taken had they seen the original image. To approximate the loss function for this model, we train a discriminator that tries to distinguish whether a user's action was taken in response to the compressed image or the original. We evaluate our method through experiments with human participants on four tasks: reading handwritten digits, verifying photos of faces, browsing an online shopping catalogue, and playing a car racing video game. The results show that our method learns to match the user's actions with and without compression at lower bitrates than baseline methods, and adapts the compression model to the user's behavior: it preserves the digit number and randomizes handwriting style in the digit reading task, preserves hats and eyeglasses while randomizing faces in the photo verification task, preserves the perceived price of an item while randomizing its color and background in the online shopping task, and preserves upcoming bends in the road in the car racing game.

[Generalized Linear Bandits with Local Differential Privacy](#)

- Yuxuan Han, Zhipeng Liang, Yang Wang, Jiheng Zhang
- abstract@[open-review](#): Contextual bandit algorithms are useful in personalized online decision-making. However, many applications such as personalized medicine and online advertising require the utilization of individual-specific information for effective learning, while user's data should remain private from the server due to privacy concerns. This motivates the introduction of local differential privacy (LDP), a stringent notion in privacy, to contextual bandits. In this paper, we design LDP algorithms for stochastic generalized linear bandits to achieve the same regret bound as in non-privacy settings. Our main idea is to develop a stochastic gradient-based estimator and update mechanism to ensure LDP. We then exploit the flexibility of stochastic gradient descent (SGD), whose theoretical guarantee for bandit problems is rarely explored, in dealing with generalized linear bandits. We also develop an estimator and update mechanism based on Ordinary Least Square (OLS) for linear bandits. Finally, we conduct experiments with both simulation and real-world datasets to demonstrate the consistently superb performance of our algorithms under LDP constraints with reasonably small parameters $\$(\backslash w_{\text{epsilon}}, \backslash delta)$ to ensure strong privacy protection.

[On the Algorithmic Stability of Adversarial Training](#)

- Yue Xing, Qifan Song, Guang Cheng
- abstract@[open-review](#): The adversarial training is a popular tool to remedy the vulnerability of deep learning models against adversarial attacks, and there is rich theoretical literature on the training loss of adversarial training algorithms. In contrast, this paper studies the algorithmic stability of a generic

adversarial training algorithm, which can further help to establish an upper bound for generalization error. By figuring out the stability upper bound and lower bound, we argue that the non-differentiability issue of adversarial training causes worse algorithmic stability than their natural counterparts. To tackle this problem, we consider a noise injection method. While the non-differentiability problem seriously affects the stability of adversarial training, injecting noise enables the training trajectory to avoid the occurrence of non-differentiability with dominating probability, hence enhancing the stability performance of adversarial training. Our analysis also studies the relation between the algorithm stability and numerical approximation error of adversarial attacks.

[Width-based Lookaheads with Learnt Base Policies and Heuristics Over the Atari-2600 Benchmark](#)

- Stefan O'Toole, Nir Lipovetzky, Miquel Ramirez, Adrian Pearce
- abstract@[open-review](#): We propose new width-based planning and learning algorithms inspired from a careful analysis of the design decisions made by previous width-based planners. The algorithms are applied over the Atari-2600 games and our best performing algorithm, Novelty guided Critical Path Learning (N-CPL), outperforms the previously introduced width-based planning and learning algorithms $\$pi\$-IW(1)$, $\$pi\$-IW(1)+$ and $\$pi\$-HIW(n, 1)$. Furthermore, we present a taxonomy of the Atari-2600 games according to some of their defining characteristics. This analysis of the games provides further insight into the behaviour and performance of the algorithms introduced. Namely, for games with large branching factors, and games with sparse meaningful rewards, N-CPL outperforms $\$pi\$-IW$, $\$pi\$-IW(1)+$ and $\$pi\$-HIW(n, 1)$.

[Characterizing possible failure modes in physics-informed neural networks](#)

- Aditi Krishnapriyan, Amir Gholami, Shandian Zhe, Robert Kirby, Michael W. Mahoney
- abstract@[open-review](#): Recent work in scientific machine learning has developed so-called physics-informed neural network (PINN) models. The typical approach is to incorporate physical domain knowledge as soft constraints on an empirical loss function and use existing machine learning methodologies to train the model. We demonstrate that, while existing PINN methodologies can learn good models for relatively trivial problems, they can easily fail to learn relevant physical phenomena for even slightly more complex problems. In particular, we analyze several distinct situations of widespread physical interest, including learning differential equations with convection, reaction, and diffusion operators. We provide evidence that the soft regularization in PINNs, which involves PDE-based differential operators, can introduce a number of subtle problems, including making the problem more ill-conditioned. Importantly, we show that these possible failure modes are not due to the lack of expressivity in the NN architecture, but that the PINN's setup makes the loss landscape very hard to optimize. We then describe two promising solutions to address these failure modes. The first approach is to use curriculum regularization, where the PINN's loss term starts from a simple PDE regularization, and becomes progressively more complex as the NN gets trained. The second approach is to pose the problem as a sequence-to-sequence learning task, rather than learning to predict the entire space-time at once. Extensive testing shows that we can achieve up to 1-2 orders of magnitude lower error with these methods as compared to regular PINN training.

[Artistic Style Transfer with Internal-external Learning and Contrastive Learning](#)

- Haibo Chen, lei zhao, Zhizhong Wang, Huiming Zhang, Zhiwen Zuo, Ailin Li, Wei Xing, Dongming Lu
- abstract@[open-review](#): Although existing artistic style transfer methods have achieved significant improvement with deep neural networks, they still suffer from artifacts such as disharmonious colors and repetitive patterns. Motivated by this, we propose an internal-external style transfer method with two contrastive losses. Specifically, we utilize internal statistics of a single style image to determine the colors and texture patterns of the stylized image, and in the meantime, we leverage the external information of the large-scale style dataset to learn the human-aware style information, which makes the color distributions and texture patterns in the stylized image more reasonable and harmonious. In addition, we argue that existing style transfer methods only consider the content-to-stylization and style-to-stylization relations, neglecting the stylization-to-stylization relations. To address this issue, we introduce two contrastive losses, which pull the multiple stylization embeddings closer to each other when they share the same content or style, but push far away otherwise. We conduct extensive experiments, showing that our proposed method can not only produce visually more harmonious and satisfying artistic images, but also promote the stability and consistency of rendered video clips.

[Fast Abductive Learning by Similarity-based Consistency Optimization](#)

- Yu-Xuan Huang, Wang-Zhou Dai, Le-Wen Cai, Stephen H Muggleton, Yuan Jiang
- abstract@[open-review](#): To utilize the raw inputs and symbolic knowledge simultaneously, some recent neuro-symbolic learning methods use abduction, i.e., abductive reasoning, to integrate sub-symbolic perception and logical inference. While the perception model, e.g., a neural network, outputs some facts that are inconsistent with the symbolic background knowledge base, abduction can help revise the incorrect perceived facts by minimizing the inconsistency between them and the background knowledge. However, to enable effective abduction, previous approaches need an initialized perception model that discriminates the input raw instances. This limits the application of these methods, as the discrimination ability is usually acquired from a thorough pre-training when the raw inputs are difficult to classify. In this paper, we propose a novel abduction strategy, which leverages the similarity between samples, rather than the output information by the perceptual neural network, to guide the search in abduction. Based on this principle, we further present ABductive Learning with Similarity (ABLSim) and apply it to some difficult neuro-symbolic learning tasks. Experiments show that the efficiency of ABLSim is significantly higher than the state-of-the-art neuro-symbolic methods, allowing it to achieve better performance with less labeled data and weaker domain knowledge.

[To Beam Or Not To Beam: That is a Question of Cooperation for Language GANs](#)

- Thomas Scialom, Paul-Alexis Dray, Jacopo Staiano, Sylvain Lamprier, Benjamin Piwowarski
- abstract@[open-review](#): Due to the discrete nature of words, language GANs require to be optimized from rewards provided by discriminator networks, via reinforcement learning methods. This is a much harder setting than for continuous tasks, which enjoy gradient flows from discriminators to generators, usually leading to dramatic learning instabilities. However, we claim that this can be solved by making discriminator and generator networks cooperate to produce output sequences during training. These cooperative outputs, inherently built to obtain higher discrimination scores, not only provide denser rewards for training but also form a more compact artificial set for discriminator training, hence improving its accuracy and stability. In this paper, we show that our SelfGAN framework, built on this cooperative principle, outperforms Teacher Forcing and obtains state-of-the-art results on two challenging tasks, Summarization and Question Generation.

[Shapley Residuals: Quantifying the limits of the Shapley value for explanations](#)

- Indra Kumar, Carlos Scheidegger, Suresh Venkatasubramanian, Sorelle Friedler
- abstract@[open-review](#): Popular feature importance techniques compute additive approximations to nonlinear models by first defining a cooperative game describing the value of different subsets of the model's features, then calculating the resulting game's Shapley values to attribute credit additively between the features. However, the specific modeling settings in which the Shapley values are a poor approximation for the true game have not been well-described. In this paper we utilize an interpretation of Shapley values as the result of an orthogonal projection between vector spaces to calculate a residual representing the kernel component of that projection. We provide an algorithm for computing these residuals, characterize different modeling settings based on the value of the residuals, and demonstrate that they capture information about model predictions that Shapley values cannot. Shapley residuals can thus act as a warning to practitioners against overestimating the degree to which Shapley-value-based explanations give them insight into a model.

[The Elastic Lottery Ticket Hypothesis](#)

- Xiaohan Chen, Yu Cheng, Shuohang Wang, Zhe Gan, Jingjing Liu, Zhangyang Wang
- abstract@[open-review](#): Lottery Ticket Hypothesis (LTH) raises keen attention to identifying sparse trainable subnetworks, or winning tickets, which can be trained in isolation to achieve similar or even better performance compared to the full models. Despite many efforts being made, the most effective method to identify such winning tickets is still Iterative Magnitude-based Pruning (IMP), which is computationally expensive and has to be run thoroughly for every different network. A natural question that comes in is: can we “transform” the winning ticket found in one network to another with a different architecture, yielding a winning ticket for the latter at the beginning, without re-doing the expensive IMP? Answering this question is not only practically relevant for efficient “once-for-all” winning ticket finding, but also theoretically appealing for uncovering inherently scalable sparse patterns in networks. We conduct extensive experiments on CIFAR-10 and ImageNet, and propose a variety of strategies to tweak the winning tickets found from different networks of the same model family (e.g., ResNets). Based on these results, we articulate the Elastic Lottery Ticket Hypothesis (E-LTH): by mindfully replicating (or dropping) and re-ordering layers for one network, its corresponding winning ticket could be stretched (or squeezed) into a subnetwork for another deeper (or shallower) network from the same family, whose performance is nearly the same competitive as the latter’s winning ticket directly found by IMP. We have also extensively compared E-LTH with pruning-at-initialization and dynamic sparse training methods, as well as discussed the generalizability of E-LTH to different model families, layer types, and across datasets. Code is available at <https://github.com/VITA-Group/ElasticLTH>.

[Joint Inference for Neural Network Depth and Dropout Regularization](#)

- Kishan K C, Rui Li, MohammadMahdi Gilany
- abstract@[open-review](#): Dropout regularization methods prune a neural network's pre-determined backbone structure to avoid overfitting. However, a deep model still tends to be poorly calibrated with high confidence on incorrect predictions. We propose a unified Bayesian model selection method to jointly infer the most plausible network depth warranted by data, and perform dropout regularization simultaneously. In particular, to infer network depth we define a beta process over the number of hidden layers which allows it to go to infinity. Layer-wise activation probabilities induced by the beta process modulate neuron activation via binary vectors of a conjugate Bernoulli process. Experiments across domains show that by adapting network depth and dropout regularization to data, our method achieves superior performance comparing to state-of-the-art methods with well-calibrated uncertainty estimates. In continual learning, our method enables neural networks to dynamically evolve their depths to accommodate incrementally available data beyond their initial structures, and alleviate catastrophic forgetting.

[Tractable Density Estimation on Learned Manifolds with Conformal Embedding Flows](#)

- Brendan Ross, Jesse Cresswell
- abstract@[open-review](#): Normalizing flows are generative models that provide tractable density estimation via an invertible transformation from a simple base distribution to a complex target distribution. However, this technique cannot directly model data supported on an unknown low-dimensional manifold, a common occurrence in real-world domains such as image data. Recent attempts to remedy this limitation have introduced geometric complications that defeat a central benefit of normalizing flows: exact density estimation. We recover this benefit with Conformal Embedding Flows, a framework for designing flows that learn manifolds with tractable densities. We argue that composing a standard flow with a trainable conformal embedding is the most natural way to model manifold-supported data. To this end, we present a series of conformal building blocks and apply them in experiments with synthetic and real-world data to demonstrate that flows can model manifold-supported distributions without sacrificing tractable likelihoods.

[The Limits of Optimal Pricing in the Dark](#)

- Quinlan Dawkins, Minbiao Han, Haifeng Xu
- abstract@[open-review](#): A ubiquitous learning problem in today's digital market is, during repeated interactions between a seller and a buyer, how a seller can gradually learn optimal pricing decisions based on the buyer's past purchase responses. A fundamental challenge of learning in such a strategic setup is that the buyer will naturally have incentives to manipulate his responses in order to induce more favorable learning outcomes for him. To understand the limits of the seller's learning when facing such a strategic and possibly manipulative buyer, we study a natural yet powerful buyer manipulation strategy. That is, before the pricing game starts, the buyer simply commits to “imitate” a different value function by pretending to always react optimally according to this imitative value function. We fully characterize the optimal imitative value function that the buyer should imitate as well as the resultant seller revenue and buyer surplus under this optimal buyer manipulation. Our characterizations reveal many useful insights about what happens at equilibrium. For example, a seller with concave production cost will obtain essentially 0 revenue at equilibrium whereas the revenue for a seller with convex production cost is the Bregman divergence of her cost function between no production and certain production. Finally, and importantly, we show that a more powerful class of pricing schemes does not necessarily increase, in fact, may be harmful to, the seller's revenue. Our results not only lead to an effective prescriptive way for buyers to manipulate learning algorithms but also shed lights on the limits of what a seller can really achieve when pricing in the dark.

[No RL, No Simulation: Learning to Navigate without Navigating](#)

- Meera Hahn, Devendra Singh Chaplot, Shubham Tulsiani, Mustafa Mukadam, James M. Rehg, Abhinav Gupta
- abstract@[open-review](#): Most prior methods for learning navigation policies require access to simulation environments, as they need online policy interaction and rely on ground-truth maps for rewards. However, building simulators is expensive (requires manual effort for each and every scene) and creates challenges in transferring learned policies to robotic platforms in the real-world, due to the sim-to-real domain gap. In this paper, we pose a simple question: Do we really need active interaction, ground-truth maps or even reinforcement-learning (RL) in order to solve the image-goal navigation task? We propose a self-supervised approach to learn to navigate from only passive videos of roaming. Our approach, No RL, No Simulator (NRNS), is simple and scalable, yet highly effective. NRNS outperforms RL-based formulations by a significant margin. We present NRNS as a strong baseline for any future image-based navigation tasks that use RL or Simulation.

[Analogous to Evolutionary Algorithm: Designing a Unified Sequence Model](#)

- Jiangning Zhang, Chao Xu, Jian Li, Wenzhou Chen, Yabiao Wang, Ying Tai, Shuo Chen, Chengjie Wang, Feiyue Huang, Yong Liu
- abstract@[open-review](#): Inspired by biological evolution, we explain the rationality of Vision Transformer by analogy with the proven practical Evolutionary Algorithm (EA) and derive that both of them have consistent mathematical representation. Analogous to the dynamic local population in EA, we improve the existing transformer structure and propose a more efficient EAT model, and design task-related heads to deal with different tasks more flexibly. Moreover, we introduce the spatial-filling curve into the current vision transformer to sequence image data into a uniform sequential format. Thus we can design a unified EAT framework to address multi-modal tasks, separating the network architecture from the data format adaptation. Our approach achieves state-of-the-art results on the ImageNet classification task compared with recent vision transformer works while having smaller parameters and greater throughput. We further conduct multi-modal tasks to demonstrate the superiority of the unified EAT, \eg, Text-Based Image Retrieval, and our approach improves the rank-1 by +3.7 points over the baseline on the CSS dataset.

[Improving Compositionalty of Neural Networks by Decoding Representations to Inputs](#)

- Mike Wu, Noah Goodman, Stefano Ermon
- abstract@[open-review](#): In traditional software programs, it is easy to trace program logic from variables back to input, apply assertion statements to block erroneous behavior, and compose programs together. Although deep learning programs have demonstrated strong performance on novel applications, they sacrifice many of the functionalities of traditional software programs. With this as motivation, we take a modest first step towards improving deep learning programs by jointly training a generative model to constrain neural network activations to "decode" back to inputs. We call this design a Decodable Neural Network, or DecNN. Doing so enables a form of compositionality in neural networks, where one can recursively compose DecNN with itself to create an ensemble-like model with uncertainty. In our experiments, we demonstrate applications of this uncertainty to out-of-distribution detection, adversarial example detection, and calibration --- while matching standard neural networks in accuracy. We further explore this compositionality by combining DecNN with pretrained models, where we show promising results that neural networks can be regularized from using protected features.

[The Hardness Analysis of Thompson Sampling for Combinatorial Semi-bandits with Greedy Oracle](#)

- Fang Kong, Yueran Yang, Wei Chen, Shuai Li
- abstract@[open-review](#): Thompson sampling (TS) has attracted a lot of interest in the bandit area. It was introduced in the 1930s but has not been theoretically proven until recent years. All of its analysis in the combinatorial multi-armed bandit (CMAB) setting requires an exact oracle to provide optimal solutions with any input. However, such an oracle is usually not feasible since many combinatorial optimization problems are NP-hard and only approximation oracles are available. An example \cite{WangC18} has shown the failure of TS to learn with an approximation oracle. However, this oracle is uncommon and is designed only for a specific problem instance. It is still an open question whether the convergence analysis of TS can be extended beyond the exact oracle in CMAB. In this paper, we study this question under the greedy oracle, which is a common (approximation) oracle with theoretical guarantees to solve many (offline) combinatorial optimization problems. We provide a problem-dependent regret lower bound of order $\$O(\Omega(\log T/\Delta^2))$ to quantify the hardness of TS to solve CMAB problems with greedy oracle, where T is the time horizon and Δ is some reward gap. We also provide an almost matching regret upper bound. These are the first theoretical results for TS to solve CMAB with a common approximation oracle and break the misconception that TS cannot work with approximation oracles.

[Universal Semi-Supervised Learning](#)

- Zhuo Huang, Chao Xue, Bo Han, Jian Yang, Chen Gong
- abstract@[open-review](#): Universal Semi-Supervised Learning (UniSSL) aims to solve the open-set problem where both the class distribution (i.e., class set) and feature distribution (i.e., feature domain) are different between labeled dataset and unlabeled dataset. Such a problem seriously hinders the realistic landing of classical SSL. Different from the existing SSL methods targeting at the open-set problem that only study one certain scenario of class distribution mismatch and ignore the feature distribution mismatch, we consider a more general case where a mismatch exists in both class and feature distribution. In this case, we propose a "Class-shAring data detection and Feature Adaptation" (CAFA) framework which requires no prior knowledge of the class relationship between the labeled dataset and unlabeled dataset. Particularly, CAFA utilizes a novel scoring strategy to detect the data in the shared class set. Then, it conducts domain adaptation to fully exploit the value of the detected class-sharing data for better semi-supervised consistency training. Exhaustive experiments on several benchmark datasets show the effectiveness of our method in tackling open-set problems.

[Improving Deep Learning Interpretability by Saliency Guided Training](#)

- Aya Abdelsalam Ismail, Hector Corrada Bravo, Soheil Feizi
- abstract@[open-review](#): Saliency methods have been widely used to highlight important input features in model predictions. Most existing methods use backpropagation on a modified gradient function to generate saliency maps. Thus, noisy gradients can result in unfaithful feature attributions. In this paper, we tackle this issue and introduce a \{it saliency guided training\} procedure for neural networks to reduce noisy gradients used in predictions while retaining the predictive performance of the model. Our saliency guided training procedure iteratively masks features with small and potentially noisy gradients while maximizing the similarity of model outputs for both masked and unmasked inputs. We apply the saliency guided training procedure to various synthetic and real data sets from computer vision, natural language processing, and time series across diverse neural architectures, including Recurrent Neural Networks, Convolutional Networks, and Transformers. Through qualitative and quantitative evaluations, we show that saliency guided training procedure significantly improves model interpretability across various domains while preserving its predictive performance.

[SurvITE: Learning Heterogeneous Treatment Effects from Time-to-Event Data](#)

- Alicia Curth, Changhee Lee, Mihaela van der Schaar
- abstract@[open-review](#): We study the problem of inferring heterogeneous treatment effects from time-to-event data. While both the related problems of (i) estimating treatment effects for binary or continuous outcomes and (ii) predicting survival outcomes have been well studied in the recent machine learning literature, their combination -- albeit of high practical relevance -- has received considerably less attention. With the ultimate goal of reliably estimating the effects of treatments on instantaneous risk and survival probabilities, we focus on the problem of learning (discrete-time) treatment-specific conditional hazard functions. We find that unique challenges arise in this context due to a variety of covariate shift issues that go beyond a mere combination of well-studied confounding and censoring biases. We theoretically analyse their effects by adapting recent generalization bounds from domain adaptation and treatment effect estimation to our setting and discuss implications for model design. We use the resulting insights to propose a novel deep learning method for treatment-specific hazard estimation based on balancing representations. We investigate performance across a range of experimental settings and empirically confirm that our method outperforms baselines by addressing covariate shifts from various sources.

[Optimal Rates for Nonparametric Density Estimation under Communication Constraints](#)

- Jayadev Acharya, Clement Canonne, Aditya Vikram Singh, Himanshu Tyagi
- abstract@[open-review](#): We consider density estimation for Besov spaces when the estimator is restricted to use only a limited number of bits about each sample. We provide a noninteractive adaptive estimator which exploits the sparsity of wavelet bases, along with a simulate-and-infer technique from parametric estimation under communication constraints. We show that our estimator is nearly rate-optimal by deriving minmax lower bounds that hold even when interactive protocols are allowed. Interestingly, while our wavelet-based estimator is almost rate-optimal for Sobolev spaces as well, it is unclear whether the standard Fourier basis, which arise naturally for those spaces, can be used to achieve the same performance.

[Rank Overspecified Robust Matrix Recovery: Subgradient Method and Exact Recovery](#)

- Lijun Ding, Liwei Jiang, Yudong Chen, Qing Qu, Zhihui Zhu
- abstract@[open-review](#): We study the robust recovery of a low-rank matrix from sparsely and grossly corrupted Gaussian measurements, with no prior knowledge on the intrinsic rank. We consider the robust matrix factorization approach. We employ a robust ℓ_1 loss function and deal with the challenge of the unknown rank by using an overspecified factored representation of the matrix variable. We then solve the associated nonconvex nonsmooth problem using a subgradient method with diminishing stepsizes. We show that under a regularity condition on the sensing matrices and corruption, which we call restricted direction preserving property (RDPP), even with rank overspecified, the subgradient method converges to the exact low-rank solution at a sublinear rate. Moreover, our result is more general in the sense that it automatically speeds up to a linear rate once the factor rank matches the unknown rank. On the other hand, we show that the RDPP condition holds under generic settings, such as Gaussian measurements under independent or adversarial sparse corruptions, where the result could be of independent interest. Both the exact recovery and the convergence rate of the proposed subgradient method are numerically verified in the overspecified regime. Moreover, our experiment further shows that our particular design of

diminishing stepsize effectively prevents overfitting for robust recovery under overparameterized models, such as robust matrix sensing and learning robust deep image prior. This regularization effect is worth further investigation.

[Improving Computational Efficiency in Visual Reinforcement Learning via Stored Embeddings](#)

- Lili Chen, Kimin Lee, Aravind Srinivas, Pieter Abbeel
- abstract@[open-review](#): Recent advances in off-policy deep reinforcement learning (RL) have led to impressive success in complex tasks from visual observations. Experience replay improves sample-efficiency by reusing experiences from the past, and convolutional neural networks (CNNs) process high-dimensional inputs effectively. However, such techniques demand high memory and computational bandwidth. In this paper, we present Stored Embeddings for Efficient Reinforcement Learning (SEER), a simple modification of existing off-policy RL methods, to address these computational and memory requirements. To reduce the computational overhead of gradient updates in CNNs, we freeze the lower layers of CNN encoders early in training due to early convergence of their parameters. Additionally, we reduce memory requirements by storing the low-dimensional latent vectors for experience replay instead of high-dimensional images, enabling an adaptive increase in the replay buffer capacity, a useful technique in constrained-memory settings. In our experiments, we show that SEER does not degrade the performance of RL agents while significantly saving computation and memory across a diverse set of DeepMind Control environments and Atari games.

[Learning Generalized Gumbel-max Causal Mechanisms](#)

- Guy Lorberbom, Daniel D. Johnson, Chris J. Maddison, Daniel Tarlow, Tamir Hazan
- abstract@[open-review](#): To perform counterfactual reasoning in Structural Causal Models (SCMs), one needs to know the causal mechanisms, which provide factorizations of conditional distributions into noise sources and deterministic functions mapping realizations of noise to samples. Unfortunately, the causal mechanism is not uniquely identified by data that can be gathered by observing and interacting with the world, so there remains the question of how to choose causal mechanisms. In recent work, Oberst & Sontag (2019) propose Gumbel-max SCMs, which use Gumbel-max reparameterizations as the causal mechanism due to an appealing counterfactual stability property. However, the justification requires appealing to intuition. In this work, we instead argue for choosing a causal mechanism that is best under a quantitative criteria such as minimizing variance when estimating counterfactual treatment effects. We propose a parameterized family of causal mechanisms that generalize Gumbel-max. We show that they can be trained to minimize counterfactual effect variance and other losses on a distribution of queries of interest, yielding lower variance estimates of counterfactual treatment effect than fixed alternatives, also generalizing to queries not seen at training time.

[Bandit Learning with Delayed Impact of Actions](#)

- Wei Tang, Chien-Ju Ho, Yang Liu
- abstract@[open-review](#): We consider a stochastic multi-armed bandit (MAB) problem with delayed impact of actions. In our setting, actions taken in the past impact the arm rewards in the subsequent future. This delayed impact of actions is prevalent in the real world. For example, the capability to pay back a loan for people in a certain social group might depend on historically how frequently that group has been approved loan applications. If banks keep rejecting loan applications to people in a disadvantaged group, it could create a feedback loop and further damage the chance of getting loans for people in that group. In this paper, we formulate this delayed and long-term impact of actions within the context of multi-armed bandits. We generalize the bandit setting to encode the dependency of this ``bias'' due to the action history during learning. The goal is to maximize the collected utilities over time while taking into account the dynamics created by the delayed impacts of historical actions. We propose an algorithm that achieves a regret of $\tilde{O}(KT^{2/3})$ and show a matching regret lower bound of $\Omega(KT^{2/3})$, where K is the number of arms and T is the learning horizon. Our results complement the bandit literature by adding techniques to deal with actions with long-term impacts and have implications in designing fair algorithms.

[A Stochastic Newton Algorithm for Distributed Convex Optimization](#)

- Brian Bullins, Kshitij Patel, Ohad Shamir, Nathan Srebro, Blake E. Woodworth
- abstract@[open-review](#): We propose and analyze a stochastic Newton algorithm for homogeneous distributed stochastic convex optimization, where each machine can calculate stochastic gradients of the same population objective, as well as stochastic Hessian-vector products (products of an independent unbiased estimator of the Hessian of the population objective with arbitrary vectors), with many such stochastic computations performed between rounds of communication. We show that our method can reduce the number, and frequency, of required communication rounds, compared to existing methods without hurting performance, by proving convergence guarantees for quasi-self-concordant objectives (e.g., logistic regression), alongside empirical evidence.

[Are Transformers more robust than CNNs?](#)

- Yutong Bai, Jieru Mei, Alan L. Yuille, Cihang Xie
- abstract@[open-review](#): Transformer emerges as a powerful tool for visual recognition. In addition to demonstrating competitive performance on a broad range of visual benchmarks, recent works also argue that Transformers are much more robust than Convolutional Neural Networks (CNNs). Nonetheless, surprisingly, we find these conclusions are drawn from unfair experimental settings, where Transformers and CNNs are compared at different scales and are applied with distinct training frameworks. In this paper, we aim to provide the first fair & in-depth comparisons between Transformers and CNNs, focusing on robustness evaluations. With our unified training setup, we first challenge the previous belief that Transformers outshine CNNs when measuring adversarial robustness. More surprisingly, we find CNNs can easily be as robust as Transformers on defending against adversarial attacks, if they properly adopt Transformers' training recipes. While regarding generalization on out-of-distribution samples, we show pre-training on (external) large-scale datasets is not a fundamental request for enabling Transformers to achieve better performance than CNNs. Moreover, our ablations suggest such stronger generalization is largely benefited by the Transformer's self-attention-like architectures per se, rather than by other training setups. We hope this work can help the community better understand and benchmark the robustness of Transformers and CNNs. The code and models are publicly available at: <https://github.com/ytongbai/ViTvs-CNNs>.

[Towards Sharper Generalization Bounds for Structured Prediction](#)

- Shaojie Li, Yong Liu
- abstract@[open-review](#): In this paper, we investigate the generalization performance of structured prediction learning and obtain state-of-the-art generalization bounds. Our analysis is based on factor graph decomposition of structured prediction algorithms, and we present novel margin guarantees from three different perspectives: Lipschitz continuity, smoothness, and space capacity condition. In the Lipschitz continuity scenario, we improve the square-root dependency on the label set cardinality of existing bounds to a logarithmic dependence. In the smoothness scenario, we provide generalization bounds that are not only a logarithmic dependency on the label set cardinality but a faster convergence rate of order $\mathcal{O}(\frac{1}{n})$ on the sample size n . In the space capacity scenario, we obtain bounds that do not depend on the label set cardinality and have faster convergence rates than $\mathcal{O}(\frac{1}{\sqrt{n}})$. In each scenario, applications are provided to suggest that these conditions are easy to be satisfied.

[Automated Discovery of Adaptive Attacks on Adversarial Defenses](#)

- Chengyuan Yao, Pavol Bielik, Petar Tsankov, Martin Vechev
- abstract@[open-review](#): Reliable evaluation of adversarial defenses is a challenging task, currently limited to an expert who manually crafts attacks that exploit the defense's inner workings, or to approaches based on ensemble of fixed attacks, none of which may be effective for the specific defense at hand. Our key observation is that adaptive attacks are composed from a set of reusable building blocks that can be formalized in a search space and used to automatically discover attacks for unknown defenses. We evaluated our approach on 24 adversarial defenses and show that it outperforms AutoAttack, the current state-of-the-art tool for reliable evaluation of adversarial defenses: our tool discovered significantly stronger attacks by producing 3.0%-50.8% additional adversarial examples for 10 models, while obtaining attacks with slightly stronger or similar strength for the remaining models.

[PolarStream: Streaming Object Detection and Segmentation with Polar Pillars](#)

- Qi Chen, Sourabh Vora, Oscar Beijbom
- abstract@[open-review](#): Recent works recognized lidars as an inherently streaming data source and showed that the end-to-end latency of lidar perception models can be reduced significantly by operating on wedge-shaped point cloud sectors rather than the full point cloud. However, due to use of cartesian coordinate systems these methods represent the sectors as rectangular regions, wasting memory and compute. In this work we propose using a polar coordinate system and make two key improvements on this design. First, we increase the spatial context by using multi-scale padding from neighboring sectors: preceding sector from the current scan and/or the following sector from the past scan. Second, we improve the core polar convolutional architecture by introducing feature undistortion and range stratified convolutions. Experimental results on the nuScenes dataset show significant improvements over other streaming based methods. We also achieve comparable results to existing non-streaming methods but with lower latencies.

[Representation Costs of Linear Neural Networks: Analysis and Design](#)

- Zhen Dai, Mina Karzand, Nathan Srebro
- abstract@[open-review](#): For different parameterizations (mappings from parameters to predictors), we study the regularization cost in predictor space induced by $\| \cdot \|_2$ regularization on the parameters (weights). We focus on linear neural networks as parameterizations of linear predictors. We identify the representation cost of certain sparse linear ConvNets and residual networks. In order to get a better understanding of how the architecture and parameterization affect the representation cost, we also study the reverse problem, identifying which regularizers on linear predictors (e.g., $\| \cdot \|_p$ norms, group norms, the $\| \cdot \|_{k\text{-support}}$ -norm, elastic net) can be the representation cost induced by simple $\| \cdot \|_2$ regularization, and designing the parameterizations that do so.

[Teaching via Best-Case Counterexamples in the Learning-with-Equivalence-Queries Paradigm](#)

- Akash Kumar, Yuxin Chen, Adish Singla
- abstract@[open-review](#): We study the sample complexity of teaching, termed as "teaching dimension" (TD) in the literature, for the learning-with-equivalence-queries (LwEQ) paradigm. More concretely, we consider a learner who asks equivalence queries (i.e., "is the queried hypothesis the target hypothesis?"), and a teacher responds either "yes" or "no" along with a counterexample to the queried hypothesis. This learning paradigm has been extensively studied when the learner receives worst-case or random counterexamples; in this paper, we consider the optimal teacher who picks best-case counterexamples to teach the target hypothesis within a hypothesis class. For this optimal teacher, we introduce LwEQ-TD, a notion of TD capturing the teaching complexity (i.e., the number of queries made) in this paradigm. We show that a significant reduction in queries can be achieved with best-case counterexamples, in contrast to worst-case or random counterexamples, for different hypothesis classes. Furthermore, we establish new connections of LwEQ-TD to the well-studied notions of TD in the learning-from-samples paradigm.

[Distilling Meta Knowledge on Heterogeneous Graph for Illicit Drug Trafficker Detection on Social Media](#)

- Yiyue Qian, Yiming Zhang, Yanfang (Fa) Ye, Chuxu Zhang
- abstract@[open-review](#): Driven by the considerable profits, the crime of drug trafficking (a.k.a. illicit drug trading) has co-evolved with modern technologies, e.g., social media such as Instagram has become a popular platform for marketing and selling illicit drugs. The activities of online drug trafficking are nimble and resilient, which call for novel techniques to effectively detect, disrupt, and dismantle illicit drug trades. In this paper, we propose a holistic framework named MetaHG to automatically detect illicit drug traffickers on social media (i.e., Instagram), by tackling the following two new challenges: (1) different from existing works which merely focus on analyzing post content, MetaHG is capable of jointly modeling multi-modal content and relational structured information on social media for illicit drug trafficker detection; (2) in addition, through the proposed meta-learning technique, MetaHG addresses the issue of requiring sufficient data for model training. More specifically, in our proposed MetaHG, we first build a heterogeneous graph (HG) to comprehensively characterize the complex ecosystem of drug trafficking on social media. Then, we employ a relation-based graph convolutional neural network to learn node (i.e., user) representations over the built HG, in which we introduce graph structure refinement to compensate the sparse connection among entities in the HG for more robust node representation learning. Afterwards, we propose a meta-learning algorithm for model optimization. A self-supervised module and a knowledge distillation module are further designed to exploit unlabeled data for improving the model. Extensive experiments based on the real-world data collected from Instagram demonstrate that the proposed MetaHG outperforms state-of-the-art methods.

[Curriculum Disentangled Recommendation with Noisy Multi-feedback](#)

- Hong Chen, Yudong Chen, Xin Wang, Ruobing Xie, Rui Wang, Feng Xia, Wenwu Zhu
- abstract@[open-review](#): Learning disentangled representations for user intentions from multi-feedback (i.e., positive and negative feedback) can enhance the accuracy and explainability of recommendation algorithms. However, learning such disentangled representations from multi-feedback data is challenging because i) multi-feedback is complex: there exist complex relations among different types of feedback (e.g., click, unclick, and dislike, etc) as well as various user intentions, and ii) multi-feedback is noisy: there exists noisy (useless) information both in features and labels, which may deteriorate the recommendation performance. Existing works on disentangled representation learning only focus on positive feedback, failing to handle the complex relations and noise hidden in multi-feedback data. To solve this problem, in this work we propose a Curriculum Disentangled Recommendation (CDR) model that is capable of efficiently learning disentangled representations from complex and noisy multi-feedback for better recommendation. Concretely, we design a co-filtering dynamic routing mechanism that simultaneously captures the complex relations among different behavioral feedback and user intentions as well as denoise the representations in the feature level. We then present an adjustable self-evaluating curriculum that is able to evaluate sample difficulties for better model training and conduct denoising in the label level via disregarding useless information. Our extensive experiments on several real-world datasets demonstrate that the proposed CDR model can significantly outperform several state-of-the-art methods in terms of recommendation accuracy.

[Interpretable agent communication from scratch \(with a generic visual processor emerging on the side\)](#)

- Roberto Dessì, Eugene Kharitonov, Baroni Marco
- abstract@[open-review](#): As deep networks begin to be deployed as autonomous agents, the issue of how they can communicate with each other becomes important. Here, we train two deep nets from scratch to perform realistic referent identification through unsupervised emergent communication. We show that the largely interpretable emergent protocol allows the nets to successfully communicate even about object types they did not see at training time. The visual representations induced as a by-product of our training regime, moreover, show comparable quality, when re-used as generic visual features, to a

recent self-supervised learning model. Our results provide concrete evidence of the viability of (interpretable) emergent deep net communication in a more realistic scenario than previously considered, as well as establishing an intriguing link between this field and self-supervised visual learning.

[MAU: A Motion-Aware Unit for Video Prediction and Beyond](#)

- Zheng Chang, Xinfeng Zhang, Shanshe Wang, Siwei Ma, Yan Ye, Xiang Xinguang, Wen Gao
- abstract@[open-review](#): Accurately predicting inter-frame motion information plays a key role in video prediction tasks. In this paper, we propose a Motion-Aware Unit (MAU) to capture reliable inter-frame motion information by broadening the temporal receptive field of the predictive units. The MAU consists of two modules, the attention module and the fusion module. The attention module aims to learn an attention map based on the correlations between the current spatial state and the historical spatial states. Based on the learned attention map, the historical temporal states are aggregated to an augmented motion information (AMI). In this way, the predictive unit can perceive more temporal dynamics from a wider receptive field. Then, the fusion module is utilized to further aggregate the augmented motion information (AMI) and current appearance information (current spatial state) to the final predicted frame. The computation load of MAU is relatively low and the proposed unit can be easily applied to other predictive models. Moreover, an information recalling scheme is employed into the encoders and decoders to help preserve the visual details of the predictions. We evaluate the MAU on both video prediction and early action recognition tasks. Experimental results show that the MAU outperforms the state-of-the-art methods on both tasks.

[Successor Feature Landmarks for Long-Horizon Goal-Conditioned Reinforcement Learning](#)

- Christopher Hoang, Sungryull Sohn, Jongwook Choi, Wilka Carvalho, Honglak Lee
- abstract@[open-review](#): Operating in the real-world often requires agents to learn about a complex environment and apply this understanding to achieve a breadth of goals. This problem, known as goal-conditioned reinforcement learning (GCRL), becomes especially challenging for long-horizon goals. Current methods have tackled this problem by augmenting goal-conditioned policies with graph-based planning algorithms. However, they struggle to scale to large, high-dimensional state spaces and assume access to exploration mechanisms for efficiently collecting training data. In this work, we introduce Successor Feature Landmarks (SFL), a framework for exploring large, high-dimensional environments so as to obtain a policy that is proficient for any goal. SFL leverages the ability of successor features (SF) to capture transition dynamics, using it to drive exploration by estimating state-novelty and to enable high-level planning by abstracting the state-space as a non-parametric landmark-based graph. We further exploit SF to directly compute a goal-conditioned policy for inter-landmark traversal, which we use to execute plans to "frontier" landmarks at the edge of the explored state space. We show in our experiments on MiniGrid and ViZDoom that SFL enables efficient exploration of large, high-dimensional state spaces and outperforms state-of-the-art baselines on long-horizon GCRL tasks.

[Streaming Belief Propagation for Community Detection](#)

- Yuchen Wu, Jakab Tardos, Mohammadhosseini Bateni, André Linhares, Filipe Miguel Goncalves de Almeida, Andrea Montanari, Ashkan Norouzi-Fard
- abstract@[open-review](#): The community detection problem requires to cluster the nodes of a network into a small number of well-connected ‘communities’. There has been substantial recent progress in characterizing the fundamental statistical limits of community detection under simple stochastic block models. However, in real-world applications, the network structure is typically dynamic, with nodes that join over time. In this setting, we would like a detection algorithm to perform only a limited number of updates at each node arrival. While standard voting approaches satisfy this constraint, it is unclear whether they exploit the network information optimally. We introduce a simple model for networks growing over time which we refer to as streaming stochastic block model (StSBM). Within this model, we prove that voting algorithms have fundamental limitations. We also develop a streaming belief-propagation (STREAMBP) approach, for which we prove optimality in certain regimes. We validate our theoretical findings on synthetic and real data

[The staircase property: How hierarchical structure can guide deep learning](#)

- Emmanuel Abbe, Enric Boix-Adsera, Matthew S Brennan, Guy Bresler, Dheeraj Nagaraj
- abstract@[open-review](#): This paper identifies a structural property of data distributions that enables deep neural networks to learn hierarchically. We define the “staircase” property for functions over the Boolean hypercube, which posits that high-order Fourier coefficients are reachable from lower-order Fourier coefficients along increasing chains. We prove that functions satisfying this property can be learned in polynomial time using layerwise stochastic coordinate descent on regular neural networks -- a class of network architectures and initializations that have homogeneity properties. Our analysis shows that for such staircase functions and neural networks, the gradient-based algorithm learns high-level features by greedily combining lower-level features along the depth of the network. We further back our theoretical results with experiments showing that staircase functions are learnable by more standard ResNet architectures with stochastic gradient descent. Both the theoretical and experimental results support the fact that the staircase property has a role to play in understanding the capabilities of gradient-based learning on regular networks, in contrast to general polynomial-size networks that can emulate any Statistical Query or PAC algorithm, as recently shown.

[MagNet: A Neural Network for Directed Graphs](#)

- Xitong Zhang, Yixuan He, Nathan Brugnone, Michael Perlmutter, Matthew Hirn
- abstract@[open-review](#): The prevalence of graph-based data has spurred the rapid development of graph neural networks (GNNs) and related machine learning algorithms. Yet, despite the many datasets naturally modeled as directed graphs, including citation, website, and traffic networks, the vast majority of this research focuses on undirected graphs. In this paper, we propose MagNet, a GNN for directed graphs based on a complex Hermitian matrix known as the magnetic Laplacian. This matrix encodes undirected geometric structure in the magnitude of its entries and directional information in their phase. A charge parameter attunes spectral information to variation among directed cycles. We apply our network to a variety of directed graph node classification and link prediction tasks showing that MagNet performs well on all tasks and that its performance exceeds all other methods on a majority of such tasks. The underlying principles of MagNet are such that it can be adapted to other GNN architectures.

[Hardware-adaptive Efficient Latency Prediction for NAS via Meta-Learning](#)

- Hayeon Lee, Sewoong Lee, Song Chong, Sung Ju Hwang
- abstract@[open-review](#): For deployment, neural architecture search should be hardware-aware, in order to satisfy the device-specific constraints (e.g., memory usage, latency and energy consumption) and enhance the model efficiency. Existing methods on hardware-aware NAS collect a large number of samples (e.g., accuracy and latency) from a target device, either builds a lookup table or a latency estimator. However, such approach is impractical in real-world scenarios as there exist numerous devices with different hardware specifications, and collecting samples from such a large number of devices will require prohibitive computational and monetary cost. To overcome such limitations, we propose Hardware-adaptive Efficient Latency Predictor (HELP), which formulates the device-specific latency estimation problem as a meta-learning problem, such that we can estimate the latency of a model’s performance for a given task on an unseen device with a few samples. To this end, we introduce novel hardware embeddings to embed any devices considering them as black-box functions that output latencies, and meta-learn the hardware-adaptive latency predictor in a device-dependent manner, using the hardware embeddings. We validate the proposed HELP for its latency estimation performance on unseen platforms, on which it achieves high estimation performance with as few as 10 measurement samples, outperforming all relevant baselines. We also validate end-to-end NAS frameworks using HELP against ones without it, and show that it largely reduces the total time cost of the base NAS method, in latency-constrained settings.

Topological Relational Learning on Graphs

- Yuzhou Chen, Baris Coskunuzer, Yulia Gel
- abstract@[open-review](#): Graph neural networks (GNNs) have emerged as a powerful tool for graph classification and representation learning. However, GNNs tend to suffer from over-smoothing problems and are vulnerable to graph perturbations. To address these challenges, we propose a novel topological neural framework of topological relational inference (TRI) which allows for integrating higher-order graph information to GNNs and for systematically learning a local graph structure. The key idea is to rewire the original graph by using the persistent homology of the small neighborhoods of the nodes and then to incorporate the extracted topological summaries as the side information into the local algorithm. As a result, the new framework enables us to harness both the conventional information on the graph structure and information on higher order topological properties of the graph. We derive theoretical properties on stability of the new local topological representation of the graph and discuss its implications on the graph algebraic connectivity. The experimental results on node classification tasks demonstrate that the new TRI-GNN outperforms all 14 state-of-the-art baselines on 6 out 7 graphs and exhibit higher robustness to perturbations, yielding up to 10\% better performance under noisy scenarios.

Learning Theory Can (Sometimes) Explain Generalisation in Graph Neural Networks

- Pascal Esser, Leena Chennuru Vankadara, Debarghya Ghoshdastidar
- abstract@[open-review](#): In recent years, several results in the supervised learning setting suggested that classical statistical learning-theoretic measures, such as VC dimension, do not adequately explain the performance of deep learning models which prompted a slew of work in the infinite-width and iteration regimes. However, there is little theoretical explanation for the success of neural networks beyond the supervised setting. In this paper we argue that, under some distributional assumptions, classical learning-theoretic measures can sufficiently explain generalization for graph neural networks in the transductive setting. In particular, we provide a rigorous analysis of the performance of neural networks in the context of transductive inference, specifically by analysing the generalisation properties of graph convolutional networks for the problem of node classification. While VC-dimension does result in trivial generalisation error bounds in this setting as well, we show that transductive Rademacher complexity can explain the generalisation properties of graph convolutional networks for stochastic block models. We further use the generalisation error bounds based on transductive Rademacher complexity to demonstrate the role of graph convolutions and network architectures in achieving smaller generalisation error and provide insights into when the graph structure can help in learning. The findings of this paper could re-new the interest in studying generalisation in neural networks in terms of learning-theoretic measures, albeit in specific problems.

Federated Linear Contextual Bandits

- Ruiquan Huang, Weiqiang Wu, Jing Yang, Cong Shen
- abstract@[open-review](#): This paper presents a novel federated linear contextual bandits model, where individual clients face different \$K\$-armed stochastic bandits coupled through common global parameters. By leveraging the geometric structure of the linear rewards, a collaborative algorithm called Fed-PE is proposed to cope with the heterogeneity across clients without exchanging local feature vectors or raw data. Fed-PE relies on a novel multi-client G-optimal design, and achieves near-optimal regrets for both disjoint and shared parameter cases with logarithmic communication costs. In addition, a new concept called collinearly-dependent policies is introduced, based on which a tight minimax regret lower bound for the disjoint parameter case is derived. Experiments demonstrate the effectiveness of the proposed algorithms on both synthetic and real-world datasets.

Least Square Calibration for Peer Reviews

- Sijun Tan, Jibang Wu, Xiaohui Bei, Haifeng Xu
- abstract@[open-review](#): Peer review systems such as conference paper review often suffer from the issue of miscalibration. Previous works on peer review calibration usually only use the ordinal information or assume simplistic reviewer scoring functions such as linear functions. In practice, applications like academic conferences often rely on manual methods, such as open discussions, to mitigate miscalibration. It remains an important question to develop algorithms that can handle different types of miscalibrations based on available prior knowledge. In this paper, we propose a flexible framework, namely \text{least square calibration} (LSC), for selecting top candidates from peer ratings. Our framework provably performs perfect calibration from noiseless linear scoring functions under mild assumptions, yet also provides competitive calibration results when the scoring function is from broader classes beyond linear functions and with arbitrary noise. On our synthetic dataset, we empirically demonstrate that our algorithm consistently outperforms the baseline which select top papers based on the highest average ratings.

Scaling Up Exact Neural Network Compression by ReLU Stability

- Thiago Serra, Xin Yu, Abhinav Kumar, Srikumar Ramalingam
- abstract@[open-review](#): We can compress a rectifier network while exactly preserving its underlying functionality with respect to a given input domain if some of its neurons are stable. However, current approaches to determine the stability of neurons with Rectified Linear Unit (ReLU) activations require solving or finding a good approximation to multiple discrete optimization problems. In this work, we introduce an algorithm based on solving a single optimization problem to identify all stable neurons. Our approach is on median 183 times faster than the state-of-art method on CIFAR-10, which allows us to explore exact compression on deeper (5 x 100) and wider (2 x 800) networks within minutes. For classifiers trained under an amount of L1 regularization that does not worsen accuracy, we can remove up to 56% of the connections on the CIFAR-10 dataset. The code is available at the following link, <https://github.com/yuxwind/ExactCompression>.

Passive attention in artificial neural networks predicts human visual selectivity

- Thomas Langlois, Haicheng Zhao, Erin Grant, Ishita Dasgupta, Tom Griffiths, Nori Jacoby
- abstract@[open-review](#): Developments in machine learning interpretability techniques over the past decade have provided new tools to observe the image regions that are most informative for classification and localization in artificial neural networks (ANNs). Are the same regions similarly informative to human observers? Using data from 79 new experiments and 7,810 participants, we show that passive attention techniques reveal a significant overlap with human visual selectivity estimates derived from 6 distinct behavioral tasks including visual discrimination, spatial localization, recognizability, free-viewing, cued-object search, and saliency search fixations. We find that input visualizations derived from relatively simple ANN architectures probed using guided backpropagation methods are the best predictors of a shared component in the joint variability of the human measures. We validate these correlational results with causal manipulations using recognition experiments. We show that images masked with ANN attention maps were easier for humans to classify than control masks in a speeded recognition experiment. Similarly, we find that recognition performance in the same ANN models was likewise influenced by masking input images using human visual selectivity maps. This work contributes a new approach to evaluating the biological and psychological validity of leading ANNs as models of human vision: by examining their similarities and differences in terms of their visual selectivity to the information contained in images.

GRIN: Generative Relation and Intention Network for Multi-agent Trajectory Prediction

- Longyuan Li, Jian Yao, Li Wenliang, Tong He, Tianjun Xiao, Junchi Yan, David Wipf, Zheng Zhang
- abstract@[open-review](#): Learning the distribution of future trajectories conditioned on the past is a crucial problem for understanding multi-agent systems. This is challenging because humans make decisions based on complex social relations and personal intents, resulting in highly complex uncertainties over

trajectories. To address this problem, we propose a conditional deep generative model that combines advances in graph neural networks. The prior and recognition model encodes two types of latent codes for each agent: an inter-agent latent code to represent social relations and an intra-agent latent code to represent agent intentions. The decoder is carefully devised to leverage the codes in a disentangled way to predict multi-modal future trajectory distribution. Specifically, a graph attention network built upon inter-agent latent code is used to learn continuous pair-wise relations, and an agent's motion is controlled by its latent intents and its observations of all other agents. Through experiments on both synthetic and real-world datasets, we show that our model outperforms previous work in multiple performance metrics. We also show that our model generates realistic multi-modal trajectories.

[Instance-Dependent Partial Label Learning](#)

- Ning Xu, Congyu Qiao, Xin Geng, Min-Ling Zhang
- abstract@[open-review](#): Partial label learning (PLL) is a typical weakly supervised learning problem, where each training example is associated with a set of candidate labels among which only one is true. Most existing PLL approaches assume that the incorrect labels in each training example are randomly picked as the candidate labels. However, this assumption is not realistic since the candidate labels are always instance-dependent. In this paper, we consider instance-dependent PLL and assume that each example is associated with a latent label distribution constituted by the real number of each label, representing the degree to each label describing the feature. The incorrect label with a high degree is more likely to be annotated as the candidate label. Therefore, the latent label distribution is the essential labeling information in partially labeled examples and worth being leveraged for predictive model training. Motivated by this consideration, we propose a novel PLL method that recovers the label distribution as a label enhancement (LE) process and trains the predictive model iteratively in every epoch. Specifically, we assume the true posterior density of the latent label distribution takes on the variational approximate Dirichlet density parameterized by an inference model. Then the evidence lower bound is deduced for optimizing the inference model and the label distributions generated from the variational posterior are utilized for training the predictive model. Experiments on benchmark and real-world datasets validate the effectiveness of the proposed method. Source code is available at <https://github.com/palm-ml/valen>.

[Deep Learning with Label Differential Privacy](#)

- Badih Ghazi, Noah Golowich, Ravi Kumar, Pasin Manurangsi, Chiyuan Zhang
- abstract@[open-review](#): The Randomized Response (RR) algorithm is a classical technique to improve robustness in survey aggregation, and has been widely adopted in applications with differential privacy guarantees. We propose a novel algorithm, Randomized Response with Prior (RRWithPrior), which can provide more accurate results while maintaining the same level of privacy guaranteed by RR. We then apply RRWithPrior to learn neural networks with label differential privacy (LabelDP), and show that when only the label needs to be protected, the model performance can be significantly improved over the previous state-of-the-art private baselines. Moreover, we study different ways to obtain priors, which when used with RRWithPrior can additionally improve the model performance, further reducing the accuracy gap between private and non-private models. We complement the empirical results with theoretical analysis showing that LabelDP is provably easier than protecting both the inputs and labels.

[Semialgebraic Representation of Monotone Deep Equilibrium Models and Applications to Certification](#)

- Tong Chen, Jean B. Lasserre, Victor Magron, Edouard Pauwels
- abstract@[open-review](#): Deep equilibrium models are based on implicitly defined functional relations and have shown competitive performance compared with the traditional deep networks. Monotone operator equilibrium networks (monDEQ) retain interesting performance with additional theoretical guarantees. Existing certification tools for classical deep networks cannot directly be applied to monDEQs for which much fewer tools exist. We introduce a semialgebraic representation for ReLU based monDEQs which allow to approximate the corresponding input output relation by semidefinite programs (SDP). We present several applications to network certification and obtain SDP models for the following problems : robustness certification, Lipschitz constant estimation, ellipsoidal uncertainty propagation. We use these models to certify robustness of monDEQs with respect to a general $\| \cdot \|_p$ norm. Experimental results show that the proposed models outperform existing approaches for monDEQ certification. Furthermore, our investigations suggest that monDEQs are much more robust to $\| \cdot \|_2$ perturbations than $\| \cdot \|_\infty$ perturbations.

[The Role of Global Labels in Few-Shot Classification and How to Infer Them](#)

- Ruohan Wang, Massimiliano Pontil, Carlo Ciliberto
- abstract@[open-review](#): Few-shot learning is a central problem in meta-learning, where learners must quickly adapt to new tasks given limited training data. Recently, feature pre-training has become a ubiquitous component in state-of-the-art meta-learning methods and is shown to provide significant performance improvement. However, there is limited theoretical understanding of the connection between pre-training and meta-learning. Further, pre-training requires global labels shared across tasks, which may be unavailable in practice. In this paper, we show why exploiting pre-training is theoretically advantageous for meta-learning, and in particular the critical role of global labels. This motivates us to propose Meta Label Learning (MeLa), a novel meta-learning framework that automatically infers global labels to obtain robust few-shot models. Empirically, we demonstrate that MeLa is competitive with existing methods and provide extensive ablation experiments to highlight its key properties.

[NeuS: Learning Neural Implicit Surfaces by Volume Rendering for Multi-view Reconstruction](#)

- Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, Wenping Wang
- abstract@[open-review](#): We present a novel neural surface reconstruction method, called NeuS, for reconstructing objects and scenes with high fidelity from 2D image inputs. Existing neural surface reconstruction approaches, such as DVR [Niemeyer et al., 2020] and IDR [Yariv et al., 2020], require foreground mask as supervision, easily get trapped in local minima, and therefore struggle with the reconstruction of objects with severe self-occlusion or thin structures. Meanwhile, recent neural methods for novel view synthesis, such as NeRF [Mildenhall et al., 2020] and its variants, use volume rendering to produce a neural scene representation with robustness of optimization, even for highly complex objects. However, extracting high-quality surfaces from this learned implicit representation is difficult because there are not sufficient surface constraints in the representation. In NeuS, we propose to represent a surface as the zero-level set of a signed distance function (SDF) and develop a new volume rendering method to train a neural SDF representation. We observe that the conventional volume rendering method causes inherent geometric errors (i.e. bias) for surface reconstruction, and therefore propose a new formulation that is free of bias in the first order of approximation, thus leading to more accurate surface reconstruction even without the mask supervision. Experiments on the DTU dataset and the BlendedMVS dataset show that NeuS outperforms the state-of-the-arts in high-quality surface reconstruction, especially for objects and scenes with complex structures and self-occlusion.

[Improved Guarantees for Offline Stochastic Matching via new Ordered Contention Resolution Schemes](#)

- Brian Brubach, Nathaniel Grammel, Will Ma, Aravind Srinivasan
- abstract@[open-review](#): Matching is one of the most fundamental and broadly applicable problems across many domains. In these diverse real-world applications, there is often a degree of uncertainty in the input which has led to the study of stochastic matching models. Here, each edge in the graph has a known, independent probability of existing derived from some prediction. Algorithms must probe edges to determine existence and match them irrevocably if they exist. Further, each vertex may have a patience constraint denoting how many of its neighboring edges can be probed. We present new ordered contention resolution schemes yielding improved approximation guarantees for some of the foundational problems studied in this area. For stochastic matching with patience constraints in general graphs, we provide a 0.382δ -approximate algorithm, significantly improving over the previous best 0.31δ -approximation of Baveja et al. (2018). When the vertices do not have patience constraints, we describe a 0.432δ -approximate random order probing algorithm with several corollaries such as an improved guarantee for the Prophet Secretary problem under Edge Arrivals. Finally, for the special

case of bipartite graphs with unit patience constraints on one of the partitions, we show a \$0.632\$-approximate algorithm that improves on the recent \$1/3\$-guarantee of Hikima et al. (2021).

[UFC-BERT: Unifying Multi-Modal Controls for Conditional Image Synthesis](#)

- Zhu Zhang, Jianxin Ma, Chang Zhou, Rui Men, Zhikang Li, Ming Ding, Jie Tang, Jingren Zhou, Hongxia Yang
- abstract@[open-review](#): Conditional image synthesis aims to create an image according to some multi-modal guidance in the forms of textual descriptions, reference images, and image blocks to preserve, as well as their combinations. In this paper, instead of investigating these control signals separately, we propose a new two-stage architecture, UFC-BERT, to unify any number of multi-modal controls. In UFC-BERT, both the diverse control signals and the synthesized image are uniformly represented as a sequence of discrete tokens to be processed by Transformer. Different from existing two-stage autoregressive approaches such as DALL-E and VQGAN, UFC-BERT adopts non-autoregressive generation (NAR) at the second stage to enhance the holistic consistency of the synthesized image, to support preserving specified image blocks, and to improve the synthesis speed. Further, we design a progressive algorithm that iteratively improves the non-autoregressively generated image, with the help of two estimators developed for evaluating the compliance with the controls and evaluating the fidelity of the synthesized image, respectively. Extensive experiments on a newly collected large-scale clothing dataset M2C-Fashion and a facial dataset Multi-Modal CelebA-HQ verify that UFC-BERT can synthesize high-fidelity images that comply with flexible multi-modal controls.

[Is Bang-Bang Control All You Need? Solving Continuous Control with Bernoulli Policies](#)

- Tim Seyde, Igor Gilitschenski, Wilko Schwarting, Bartolomeo Stellato, Martin Riedmiller, Markus Wulfmeier, Daniela Rus
- abstract@[open-review](#): Reinforcement learning (RL) for continuous control typically employs distributions whose support covers the entire action space. In this work, we investigate the colloquially known phenomenon that trained agents often prefer actions at the boundaries of that space. We draw theoretical connections to the emergence of bang-bang behavior in optimal control, and provide extensive empirical evaluation across a variety of recent RL algorithms. We replace the normal Gaussian by a Bernoulli distribution that solely considers the extremes along each action dimension - a bang-bang controller. Surprisingly, this achieves state-of-the-art performance on several continuous control benchmarks - in contrast to robotic hardware, where energy and maintenance cost affect controller choices. Since exploration, learning, and the final solution are entangled in RL, we provide additional imitation learning experiments to reduce the impact of exploration on our analysis. Finally, we show that our observations generalize to environments that aim to model real-world challenges and evaluate factors to mitigate the emergence of bang-bang solutions. Our findings emphasise challenges for benchmarking continuous control algorithms, particularly in light of potential real-world applications.

[Improving Generalization in Meta-RL with Imaginary Tasks from Latent Dynamics Mixture](#)

- Suyoung Lee, Sae-Young Chung
- abstract@[open-review](#): The generalization ability of most meta-reinforcement learning (meta-RL) methods is largely limited to test tasks that are sampled from the same distribution used to sample training tasks. To overcome the limitation, we propose Latent Dynamics Mixture (LDM) that trains a reinforcement learning agent with imaginary tasks generated from mixtures of learned latent dynamics. By training a policy on mixture tasks along with original training tasks, LDM allows the agent to prepare for unseen test tasks during training and prevents the agent from overfitting the training tasks. LDM significantly outperforms standard meta-RL methods in test returns on the gridworld navigation and MuJoCo tasks where we strictly separate the training task distribution and the test task distribution.

[Localization with Sampling-Argmax](#)

- Jiefeng Li, Tong Chen, Ruiqi Shi, Yujing Lou, Yong-Lu Li, Cewu Lu
- abstract@[open-review](#): Soft-argmax operation is commonly adopted in detection-based methods to localize the target position in a differentiable manner. However, training the neural network with soft-argmax makes the shape of the probability map unconstrained. Consequently, the model lacks pixel-wise supervision through the map during training, leading to performance degradation. In this work, we propose sampling-argmax, a differentiable training method that imposes implicit constraints to the shape of the probability map by minimizing the expectation of the localization error. To approximate the expectation, we introduce a continuous formulation of the output distribution and develop a differentiable sampling process. The expectation can be approximated by calculating the average error of all samples drawn from the output distribution. We show that sampling-argmax can seamlessly replace the conventional soft-argmax operation on various localization tasks. Comprehensive experiments demonstrate the effectiveness and flexibility of the proposed method. Code is available at <https://github.com/Jeff-sjtu/sampling-argmax>

[Improved Regularization and Robustness for Fine-tuning in Neural Networks](#)

- Dongyue Li, Hongyang Zhang
- abstract@[open-review](#): A widely used algorithm for transfer learning is fine-tuning, where a pre-trained model is fine-tuned on a target task with a small amount of labeled data. When the capacity of the pre-trained model is much larger than the size of the target data set, fine-tuning is prone to overfitting and "memorizing" the training labels. Hence, an important question is to regularize fine-tuning and ensure its robustness to noise. To address this question, we begin by analyzing the generalization properties of fine-tuning. We present a PAC-Bayes generalization bound that depends on the distance traveled in each layer during fine-tuning and the noise stability of the fine-tuned model. We empirically measure these quantities. Based on the analysis, we propose regularized self-labeling---the interpolation between regularization and self-labeling methods, including (i) layer-wise regularization to constrain the distance traveled in each layer; (ii) self label-correction and label-reweighting to correct mislabeled data points (that the model is confident) and reweight less confident data points. We validate our approach on an extensive collection of image and text data sets using multiple pre-trained model architectures. Our approach improves baseline methods by 1.76% (on average) for seven image classification tasks and 0.75% for a few-shot classification task. When the target data set includes noisy labels, our approach outperforms baseline methods by 3.56% on average in two noisy settings.

[BARTScore: Evaluating Generated Text as Text Generation](#)

- Weizhe Yuan, Graham Neubig, Pengfei Liu
- abstract@[open-review](#): A wide variety of NLP applications, such as machine translation, summarization, and dialog, involve text generation. One major challenge for these applications is how to evaluate whether such generated texts are actually fluent, accurate, or effective. In this work, we conceptualize the evaluation of generated text as a text generation problem, modeled using pre-trained sequence-to-sequence models. The general idea is that models trained to convert the generated text to/from a reference output or the source text will achieve higher scores when the generated text is better. We operationalize this idea using BART, an encoder-decoder based pre-trained model, and propose a metric BARTScore with a number of variants that can be flexibly applied in an unsupervised fashion to evaluation of text from different perspectives (e.g. informativeness, fluency, or factuality). BARTScore is conceptually simple and empirically effective. It can outperform existing top-scoring metrics in 16 of 22 test settings, covering evaluation of 16 datasets (e.g., machine translation, text summarization) and 7 different perspectives (e.g., informativeness, factuality). Code to calculate BARTScore is available at <https://github.com/neulab/BARTScore>, and we have released an interactive leaderboard for meta-evaluation at <http://explainaboard.nlpedia.ai/leaderboard/task-meval/> on the Explainaboard platform, which allows us to interactively understand the strengths, weaknesses, and complementarity of each metric.

[An analysis of Ermakov-Zolotukhin quadrature using kernels](#)

- Ayoub Belhadji
- abstract@[open-review](#): We study a quadrature, proposed by Ermakov and Zolotukhin in the sixties, through the lens of kernel methods. The nodes of this quadrature rule follow the distribution of a determinantal point process, while the weights are defined through a linear system, similarly to the optimal kernel quadrature. In this work, we show how these two classes of quadrature are related, and we prove a tractable formula of the expected value of the squared worst-case integration error on the unit ball of an RKHS of the former quadrature. In particular, this formula involves the eigenvalues of the corresponding kernel and leads to improving on the existing theoretical guarantees of the optimal kernel quadrature with determinantal point processes.

[Towards Understanding Why Lookahead Generalizes Better Than SGD and Beyond](#)

- Pan Zhou, Hanshu Yan, Xiaotong Yuan, Jiashi Feng, Shuicheng Yan
- abstract@[open-review](#): To train networks, lookahead algorithm $\sim\backslash$ cite{zhang2019lookahead} updates its fast weights k times via an inner-loop optimizer before updating its slow weights once by using the latest fast weights. Any optimizer, e.g. SGD, can serve as the inner-loop optimizer, and the derived lookahead generally enjoys remarkable test performance improvement over the vanilla optimizer. But theoretical understandings on the test performance improvement of lookahead remain absent yet. To solve this issue, we theoretically justify the advantages of lookahead in terms of the excess risk error which measures the test performance. Specifically, we prove that lookahead using SGD as its inner-loop optimizer can better balance the optimization error and generalization error to achieve smaller excess risk error than vanilla SGD on (strongly) convex problems and nonconvex problems with Polyak-{L}ojasiewicz condition which has been observed/proved in neural networks. Moreover, we show the stagewise optimization strategy $\sim\backslash$ cite{barshan2015stage} which decays learning rate several times during training can also benefit lookahead in improving its optimization and generalization errors on strongly convex problems. Finally, we propose a stagewise locally-regularized lookahead (SLRLA) algorithm which sums up the vanilla objective and a local regularizer to minimize at each stage and provably enjoys optimization and generalization improvement over the conventional (stagewise) lookahead. Experimental results on CIFAR10/100 and ImageNet testify its advantages. Codes is available at \url{https://github.com/sail-sg/SLRLA-optimizer}.

[Online Market Equilibrium with Application to Fair Division](#)

- Yuan Gao, Alex Peysakhovich, Christian Kroer
- abstract@[open-review](#): Computing market equilibria is a problem of both theoretical and applied interest. Much research to date focuses on the case of static Fisher markets with full information on buyers' utility functions and item supplies. Motivated by real-world markets, we consider an online setting: individuals have linear, additive utility functions; items arrive sequentially and must be allocated and priced irrevocably. We define the notion of an online market equilibrium in such a market as time-indexed allocations and prices which guarantee buyer optimality and market clearance in hindsight. We propose a simple, scalable and interpretable allocation and pricing dynamics termed as PACE. When items are drawn i.i.d. from an unknown distribution (with a possibly continuous support), we show that PACE leads to an online market equilibrium asymptotically. In particular, PACE ensures that buyers' time-averaged utilities converge to the equilibrium utilities w.r.t. a static market with item supplies being the unknown distribution and that buyers' time-averaged expenditures converge to their per-period budget. Hence, many desirable properties of market equilibrium-based fair division such as envy-freeness, Pareto optimality, and the proportional-share guarantee are also attained asymptotically in the online setting. Next, we extend the dynamics to handle quasilinear buyer utilities, which gives the first online algorithm for computing first-price pacing equilibria. Finally, numerical experiments on real and synthetic datasets show that the dynamics converges quickly under various metrics.

[Dynamic Resolution Network](#)

- Mingjian Zhu, Kai Han, Enhua Wu, Qiulin Zhang, Ying Nie, Zhenzhong Lan, Yunhe Wang
- abstract@[open-review](#): Deep convolutional neural networks (CNNs) are often of sophisticated design with numerous learnable parameters for the accuracy reason. To alleviate the expensive costs of deploying them on mobile devices, recent works have made huge efforts for excavating redundancy in pre-defined architectures. Nevertheless, the redundancy on the input resolution of modern CNNs has not been fully investigated, i.e., the resolution of input image is fixed. In this paper, we observe that the smallest resolution for accurately predicting the given image is different using the same neural network. To this end, we propose a novel dynamic-resolution network (DRNet) in which the input resolution is determined dynamically based on each input sample. Wherein, a resolution predictor with negligible computational costs is explored and optimized jointly with the desired network. Specifically, the predictor learns the smallest resolution that can retain and even exceed the original recognition accuracy for each image. During the inference, each input image will be resized to its predicted resolution for minimizing the overall computation burden. We then conduct extensive experiments on several benchmark networks and datasets. The results show that our DRNet can be embedded in any off-the-shelf network architecture to obtain a considerable reduction in computational complexity. For instance, DR-ResNet-50 achieves similar performance with an about 34% computation reduction, while gaining 1.4% accuracy increase with 10% computation reduction compared to the original ResNet-50 on ImageNet. Code will be available at <https://gitee.com/mindspore/models/tree/master/research/cv/DRNet>.

[Gauge Equivariant Transformer](#)

- Lingshen He, Yiming Dong, Yisen Wang, Dacheng Tao, Zhouchen Lin
- abstract@[open-review](#): Attention mechanism has shown great performance and efficiency in a lot of deep learning models, in which relative position encoding plays a crucial role. However, when introducing attention to manifolds, there is no canonical local coordinate system to parameterize neighborhoods. To address this issue, we propose an equivariant transformer to make our model agnostic to the orientation of local coordinate systems (\textit{i.e.}, gauge equivariant), which employs multi-head self-attention to jointly incorporate both position-based and content-based information. To enhance expressive ability, we adopt regular field of cyclic groups as feature fields in intermediate layers, and propose a novel method to parallel transport the feature vectors in these fields. In addition, we project the position vector of each point onto its local coordinate system to disentangle the orientation of the coordinate system in ambient space (\textit{i.e.}, global coordinate system), achieving rotation invariance. To the best of our knowledge, we are the first to introduce gauge equivariance to self-attention, thus name our model Gauge Equivariant Transformer (GET), which can be efficiently implemented on triangle meshes. Extensive experiments show that GET achieves state-of-the-art performance on two common recognition tasks.

[Unsupervised Object-Based Transition Models For 3D Partially Observable Environments](#)

- Antonia Creswell, Rishabh Kabra, Chris Burgess, Murray Shanahan
- abstract@[open-review](#): We present a slot-wise, object-based transition model that decomposes a scene into objects, aligns them (with respect to a slot-wise object memory) to maintain a consistent order across time, and predicts how those objects evolve over successive frames. The model is trained end-to-end without supervision using transition losses at the level of the object-structured representation rather than pixels. Thanks to the introduction of our novel alignment module, the model deals properly with two issues that are not handled satisfactorily by other transition models, namely object persistence and object identity. We show that the combination of an object-level loss and correct object alignment over time enables the model to outperform a state-of-the-art baseline, and allows it to deal well with object occlusion and re-appearance in partially observable environments.

[Robust Contrastive Learning Using Negative Samples with Diminished Semantics](#)

- Songwei Ge, Shlok Mishra, Chun-Liang Li, Haohan Wang, David Jacobs
- abstract@[open-review](#): Unsupervised learning has recently made exceptional progress because of the development of more effective contrastive learning methods. However, CNNs are prone to depend on low-level features that humans deem non-semantic. This dependency has been conjectured to induce a lack of robustness to image perturbations or domain shift. In this paper, we show that by generating carefully designed negative samples, contrastive learning can learn more robust representations with less dependence on such features. Contrastive learning utilizes positive pairs which preserve semantic information while perturbing superficial features in the training images. Similarly, we propose to generate negative samples in a reversed way, where only the superfluous instead of the semantic features are preserved. We develop two methods, texture-based and patch-based augmentations, to generate negative samples. These samples achieve better generalization, especially under out-of-domain settings. We also analyze our method and the generated texture-based samples, showing that texture features are indispensable in classifying particular ImageNet classes and especially finer classes. We also show that the model bias between texture and shape features favors them differently under different test settings.

[General Low-rank Matrix Optimization: Geometric Analysis and Sharper Bounds](#)

- Haixiang Zhang, Yingjie Bi, Javad Lavaei
- abstract@[open-review](#): This paper considers the global geometry of general low-rank minimization problems via the Burer-Monterio factorization approach. For the rank-\$1\$ case, we prove that there is no spurious second-order critical point for both symmetric and asymmetric problems if the rank-\$2\$ RIP constant δ is less than $1/2$. Combining with a counterexample with $\delta=1/2$, we show that the derived bound is the sharpest possible. For the arbitrary rank-\$r\$ case, the same property is established when the rank-\$2r\$ RIP constant δ is at most $1/3$. We design a counterexample to show that the non-existence of spurious second-order critical points may not hold if δ is at least $1/2$. In addition, for any problem with δ between $1/3$ and $1/2$, we prove that all second-order critical points have a positive correlation to the ground truth. Finally, the strict saddle property, which can lead to the polynomial-time global convergence of various algorithms, is established for both the symmetric and asymmetric problems when the rank-\$2r\$ RIP constant δ is less than $1/3$. The results of this paper significantly extend several existing bounds in the literature.

[Flow Network based Generative Models for Non-Iterative Diverse Candidate Generation](#)

- Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, Yoshua Bengio
- abstract@[open-review](#): This paper is about the problem of learning a stochastic policy for generating an object (like a molecular graph) from a sequence of actions, such that the probability of generating an object is proportional to a given positive reward for that object. Whereas standard return maximization tends to converge to a single return-maximizing sequence, there are cases where we would like to sample a diverse set of high-return solutions. These arise, for example, in black-box function optimization when few rounds are possible, each with large batches of queries, where the batches should be diverse, e.g., in the design of new molecules. One can also see this as a problem of approximately converting an energy function to a generative distribution. While MCMC methods can achieve that, they are expensive and generally only perform local exploration. Instead, training a generative policy amortizes the cost of search during training and yields to fast generation. Using insights from Temporal Difference learning, we propose GFlowNet, based on a view of the generative process as a flow network, making it possible to handle the tricky case where different trajectories can yield the same final state, e.g., there are many ways to sequentially add atoms to generate some molecular graph. We cast the set of trajectories as a flow and convert the flow consistency equations into a learning objective, akin to the casting of the Bellman equations into Temporal Difference methods. We prove that any global minimum of the proposed objectives yields a policy which samples from the desired distribution, and demonstrate the improved performance and diversity of GFlowNet on a simple domain where there are many modes to the reward function, and on a molecule synthesis task.

[Policy Finetuning: Bridging Sample-Efficient Offline and Online Reinforcement Learning](#)

- Tengyang Xie, Nan Jiang, Huan Wang, Caiming Xiong, Yu Bai
- abstract@[open-review](#): Recent theoretical work studies sample-efficient reinforcement learning (RL) extensively in two settings: learning interactively in the environment (online RL), or learning from an offline dataset (offline RL). However, existing algorithms and theories for learning near-optimal policies in these two settings are rather different and disconnected. Towards bridging this gap, this paper initiates the theoretical study of *policy finetuning*, that is, online RL where the learner has additional access to a "reference policy" μ close to the optimal policy π^* in a certain sense. We consider the policy finetuning problem in episodic Markov Decision Processes (MDPs) with S states, A actions, and horizon length H . We first design a sharp *offline reduction* algorithm---which simply executes μ and runs offline policy optimization on the collected dataset---that finds an ϵ -near-optimal policy within $\tilde{O}(H^3SC^*\epsilon^2)$ episodes, where C^* is the single-policy concentrability coefficient between μ and π^* . This offline result is the first that matches the sample complexity lower bound in this setting, and resolves a recent open question in offline RL. We then establish an $\Omega(H^3S\min\{C^*, A\}\epsilon^2)$ sample complexity lower bound for *any* policy finetuning algorithm, including those that can adaptively explore the environment. This implies that---perhaps surprisingly---the optimal policy finetuning algorithm is either offline reduction or a purely online RL algorithm that does not use μ . Finally, we design a new hybrid offline/online algorithm for policy finetuning that achieves better sample complexity than both vanilla offline reduction and purely online RL algorithms, in a relaxed setting where μ only satisfies concentrability partially up to a certain time step. Overall, our results offer a quantitative understanding on the benefit of a good reference policy, and make a step towards bridging offline and online RL.

[Reducing Information Bottleneck for Weakly Supervised Semantic Segmentation](#)

- Jungbeom Lee, Jooyoung Choi, Jisoo Mok, Sungroh Yoon
- abstract@[open-review](#): Weakly supervised semantic segmentation produces pixel-level localization from class labels; however, a classifier trained on such labels is likely to focus on a small discriminative region of the target object. We interpret this phenomenon using the information bottleneck principle: the final layer of a deep neural network, activated by the sigmoid or softmax activation functions, causes an information bottleneck, and as a result, only a subset of the task-relevant information is passed on to the output. We first support this argument through a simulated toy experiment and then propose a method to reduce the information bottleneck by removing the last activation function. In addition, we introduce a new pooling method that further encourages the transmission of information from non-discriminative regions to the classification. Our experimental evaluations demonstrate that this simple modification significantly improves the quality of localization maps on both the PASCAL VOC 2012 and MS COCO 2014 datasets, exhibiting a new state-of-the-art performance for weakly supervised semantic segmentation.

[SGD: The Role of Implicit Regularization, Batch-size and Multiple-epochs](#)

- Ayush Sekhari, Karthik Sridharan, Satyen Kale
- abstract@[open-review](#): Multi-epoch, small-batch, Stochastic Gradient Descent (SGD) has been the method of choice for learning with large over-parameterized models. A popular theory for explaining why SGD works well in practice is that the algorithm has an implicit regularization that biases its output towards a good solution. Perhaps the theoretically most well understood learning setting for SGD is that of Stochastic Convex Optimization (SCO), where it is well known that SGD learns at a rate of $O(1/\sqrt{n})$, where n is the number of samples. In this paper, we consider the problem of SCO and explore the role of implicit regularization, batch size and multiple epochs for SGD. Our main contributions are threefold: * We show that for any regularizer, there is an SCO problem for which Regularized Empirical Risk Minimization fails to learn. This automatically rules out any implicit regularization based explanation for the success of SGD. We provide a separation between SGD and learning via Gradient Descent on empirical loss (GD) in terms of sample complexity. We show that there is an SCO problem such that GD with any step size and number of iterations can only learn at a suboptimal rate: at least $\tilde{\Omega}(1/n^{5/12})$. We present a multi-epoch variant of SGD commonly used in practice. We prove that this

algorithm is at least as good as single pass SGD in the worst case. However, for certain SCO problems, taking multiple passes over the dataset can significantly outperform single pass SGD. We extend our results to the general learning setting by showing a problem which is learnable for any data distribution, and for this problem, SGD is strictly better than RERM for any regularization function. We conclude by discussing the implications of our results for deep learning, and show a separation between SGD and ERM for two layer diagonal neural networks.

[AC-GC: Lossy Activation Compression with Guaranteed Convergence](#)

- R David Evans, Tor Aamodt
- abstract@[open-review](#): Parallel hardware devices (e.g., graphics processor units) have limited high-bandwidth memory capacity. This negatively impacts the training of deep neural networks (DNNs) by increasing runtime and/or decreasing accuracy when reducing model and/or batch size to fit this capacity. Lossy compression is a promising approach to tackling memory capacity constraints, but prior approaches rely on hyperparameter search to achieve a suitable trade-off between convergence and compression, negating runtime benefits. In this paper we build upon recent developments on Stochastic Gradient Descent convergence to prove an upper bound on the expected loss increase when training with compressed activation storage. We then express activation compression error in terms of this bound, allowing the compression rate to adapt to training conditions automatically. The advantage of our approach, called AC-GC, over existing lossy compression frameworks is that, given a preset allowable increase in loss, significant compression without significant increase in error can be achieved with a single training run. When combined with error-bounded methods, AC-GC achieves 15.1x compression with an average accuracy change of 0.1% on text and image datasets. AC-GC functions on any model composed of the layers analyzed and, by avoiding compression rate search, reduces overall training time by 4.6x over SuccessiveHalving.

[Label Noise SGD Provably Prefers Flat Global Minimizers](#)

- Alex Damian, Tengyu Ma, Jason D. Lee
- abstract@[open-review](#): In overparametrized models, the noise in stochastic gradient descent (SGD) implicitly regularizes the optimization trajectory and determines which local minimum SGD converges to. Motivated by empirical studies that demonstrate that training with noisy labels improves generalization, we study the implicit regularization effect of SGD with label noise. We show that SGD with label noise converges to a stationary point of a regularized loss $\$L(\theta) + \lambda R(\theta)$, where $\$L(\theta)$ is the training loss, $\$\lambda$ is an effective regularization parameter depending on the step size, strength of the label noise, and the batch size, and $\$R(\theta)$ is an explicit regularizer that penalizes sharp minimizers. Our analysis uncovers an additional regularization effect of large learning rates beyond the linear scaling rule that penalizes large eigenvalues of the Hessian more than small ones. We also prove extensions to classification with general loss functions, significantly strengthening the prior work of Blanc et al. to global convergence and large learning rates and of HaoChen et al. to general models.

[Can we have it all? On the Trade-off between Spatial and Adversarial Robustness of Neural Networks](#)

- Sandesh Kamath, Amit Deshpande, Subrahmanyam Kambhampati Venkata, Vineeth N Balasubramanian
- abstract@[open-review](#): (Non-)robustness of neural networks to small, adversarial pixel-wise perturbations, and as more recently shown, to even random spatial transformations (e.g., translations, rotations) entreats both theoretical and empirical understanding. Spatial robustness to random translations and rotations is commonly attained via equivariant models (e.g., StdCNNs, GCNNs) and training augmentation, whereas adversarial robustness is typically achieved by adversarial training. In this paper, we prove a quantitative trade-off between spatial and adversarial robustness in a simple statistical setting. We complement this empirically by showing that: (a) as the spatial robustness of equivariant models improves by training augmentation with progressively larger transformations, their adversarial robustness worsens progressively, and (b) as the state-of-the-art robust models are adversarially trained with progressively larger pixel-wise perturbations, their spatial robustness drops progressively. Towards achieving Pareto-optimality in this trade-off, we propose a method based on curriculum learning that trains gradually on more difficult perturbations (both spatial and adversarial) to improve spatial and adversarial robustness simultaneously.

[Universal Off-Policy Evaluation](#)

- Yash Chandak, Scott Niekum, Bruno da Silva, Erik Learned-Miller, Emma Brunskill, Philip S. Thomas
- abstract@[open-review](#): When faced with sequential decision-making problems, it is often useful to be able to predict what would happen if decisions were made using a new policy. Those predictions must often be based on data collected under some previously used decision-making rule. Many previous methods enable such off-policy (or counterfactual) estimation of the expected value of a performance measure called the return. In this paper, we take the first steps towards a 'universal off-policy estimator' (UnO)---one that provides off-policy estimates and high-confidence bounds for any parameter of the return distribution. We use UnO for estimating and simultaneously bounding the mean, variance, quantiles/median, inter-quantile range, CVaR, and the entire cumulative distribution of returns. Finally, we also discuss UnO's applicability in various settings, including fully observable, partially observable (i.e., with unobserved confounders), Markovian, non-Markovian, stationary, smoothly non-stationary, and discrete distribution shifts.

[A Non-commutative Extension of Lee-Seung's Algorithm for Positive Semidefinite Factorizations](#)

- Yong Sheng Soh, Antonios Varvitsiotis
- abstract@[open-review](#): Given a data matrix $\$X \in \mathbb{R}^{m \times n}$ with non-negative entries, a Positive Semidefinite (PSD) factorization of $\$X$ is a collection of $r \times r$ -dimensional PSD matrices $\{A_{ij}\}$ and $\{B_{ij}\}$ satisfying the condition $\$X_{ij} = \text{tr}(A_{ij}^T B_{ij})$ for all $i \in [m], j \in [n]$. PSD factorizations are fundamentally linked to understanding the expressiveness of semidefinite programs as well as the power and limitations of quantum resources in information theory. The PSD factorization task generalizes the Non-negative Matrix Factorization (NMF) problem in which we seek a collection of r -dimensional non-negative vectors $\{a_i\}$ and $\{b_j\}$ satisfying $\$X_{ij} = a_i^T b_j$, for all $i \in [m], j \in [n]$ -- one can recover the latter problem by choosing matrices in the PSD factorization to be diagonal. The most widely used algorithm for computing NMFs of a matrix is the Multiplicative Update algorithm developed by Lee and Seung, in which non-negativity of the updates is preserved by scaling with positive diagonal matrices. In this paper, we describe a non-commutative extension of Lee-Seung's algorithm, which we call the Matrix Multiplicative Update (MMU) algorithm, for computing PSD factorizations. The MMU algorithm ensures that updates remain PSD by congruence scaling with the matrix geometric mean of appropriate PSD matrices, and it retains the simplicity of implementation that the multiplicative update algorithm for NMF enjoys. Building on the Majorization-Minimization framework, we show that under our update scheme the squared loss objective is non-increasing and fixed points correspond to critical points. The analysis relies on a Lieb's Concavity Theorem. Beyond PSD factorizations, we show that the MMU algorithm can be also used as a primitive to calculate block-diagonal PSD factorizations and tensor PSD factorizations. We demonstrate the utility of our method with experiments on real and synthetic data.

[Efficiently Identifying Task Groupings for Multi-Task Learning](#)

- Chris Fifty, Ehsan Amid, Zhe Zhao, Tianhe Yu, Rohan Anil, Chelsea Finn
- abstract@[open-review](#): Multi-task learning can leverage information learned by one task to benefit the training of other tasks. Despite this capacity, naively training all tasks together in one model often degrades performance, and exhaustively searching through combinations of task groupings can be prohibitively expensive. As a result, efficiently identifying the tasks that would benefit from training together remains a challenging design question without a clear solution. In this paper, we suggest an approach to select which tasks should train together in multi-task learning models. Our method determines task groupings in a single run by training all tasks together and quantifying the effect to which one task's gradient would affect another task's

loss. On the large-scale Taskonomy computer vision dataset, we find this method can decrease test loss by 10.0% compared to simply training all tasks together while operating 11.6 times faster than a state-of-the-art task grouping method.

[Instance-Conditioned GAN](#)

- Arantxa Casanova, Marlene Careil, Jakob Verbeek, Michal Drozdzal, Adriana Romero Soriano
- abstract@[open-review](#): Generative Adversarial Networks (GANs) can generate near photo realistic images in narrow domains such as human faces. Yet, modeling complex distributions of datasets such as ImageNet and COCO-Stuff remains challenging in unconditional settings. In this paper, we take inspiration from kernel density estimation techniques and introduce a non-parametric approach to modeling distributions of complex datasets. We partition the data manifold into a mixture of overlapping neighborhoods described by a datapoint and its nearest neighbors, and introduce a model, called instance-conditioned GAN (IC-GAN), which learns the distribution around each datapoint. Experimental results on ImageNet and COCO-Stuff show that IC-GAN significantly improves over unconditional models and unsupervised data partitioning baselines. Moreover, we show that IC-GAN can effortlessly transfer to datasets not seen during training by simply changing the conditioning instances, and still generate realistic images. Finally, we extend IC-GAN to the class-conditional case and show semantically controllable generation and competitive quantitative results on ImageNet; while improving over BigGAN on ImageNet-LT. Code and trained models to reproduce the reported results are available at https://github.com/facebookresearch/ic_gan.

[DeepSITH: Efficient Learning via Decomposition of What and When Across Time Scales](#)

- Brandon Jacques, Zoran Tiganj, Marc Howard, Per B Sederberg
- abstract@[open-review](#): Extracting temporal relationships over a range of scales is a hallmark of human perception and cognition---and thus it is a critical feature of machine learning applied to real-world problems. Neural networks are either plagued by the exploding/vanishing gradient problem in recurrent neural networks (RNNs) or must adjust their parameters to learn the relevant time scales (e.g., in LSTMs). This paper introduces DeepSITH, a deep network comprising biologically-inspired Scale-Invariant Temporal History (SITH) modules in series with dense connections between layers. Each SITH module is simply a set of time cells coding what happened when with a geometrically-spaced set of time lags. The dense connections between layers change the definition of what from one layer to the next. The geometric series of time lags implies that the network codes time on a logarithmic scale, enabling DeepSITH network to learn problems requiring memory over a wide range of time scales. We compare DeepSITH to LSTMs and other recent RNNs on several time series prediction and decoding tasks. DeepSITH achieves results comparable to state-of-the-art performance on these problems and continues to perform well even as the delays are increased.

[A Gaussian Process-Bayesian Bernoulli Mixture Model for Multi-Label Active Learning](#)

- Weishi Shi, Dayou Yu, Qi Yu
- abstract@[open-review](#): Multi-label classification (MLC) allows complex dependencies among labels, making it more suitable to model many real-world problems. However, data annotation for training MLC models becomes much more labor-intensive due to the correlated (hence non-exclusive) labels and a potential large and sparse label space. We propose to conduct multi-label active learning (ML-AL) through a novel integrated Gaussian Process-Bayesian Bernoulli Mixture model (GP-B²M) to accurately quantify a data sample's overall contribution to a correlated label space and choose the most informative samples for cost-effective annotation. In particular, the B²M encodes label correlations using a Bayesian Bernoulli mixture of label clusters, where each mixture component corresponds to a global pattern of label correlations. To tackle highly sparse labels under AL, the B²M is further integrated with a predictive GP to connect data features as an effective inductive bias and achieve a feature-component-label mapping. The GP predicts coefficients of mixture components that help to recover the final set of labels of a data sample. A novel auxiliary variable based variational inference algorithm is developed to tackle the non-conjugacy introduced along with the mapping process for efficient end-to-end posterior inference. The model also outputs a predictive distribution that provides both the label prediction and their correlations in the form of a label covariance matrix. A principled sampling function is designed accordingly to naturally capture both the feature uncertainty (through GP) and label covariance (through B²M) for effective data sampling. Experiments on real-world multi-label datasets demonstrate the state-of-the-art AL performance of the proposed GP-B²M model.

[Differentially Private Empirical Risk Minimization under the Fairness Lens](#)

- Cuong Tran, My Dinh, Ferdinando Fioretto
- abstract@[open-review](#): Differential Privacy (DP) is an important privacy-enhancing technology for private machine learning systems. It allows to measure and bound the risk associated with an individual participation in a computation. However, it was recently observed that DP learning systems may exacerbate bias and unfairness for different groups of individuals. This paper builds on these important observations and sheds light on the causes of the disparate impacts arising in the problem of differentially private empirical risk minimization. It focuses on the accuracy disparity arising among groups of individuals in two well-studied DP learning methods: output perturbation and differentially private stochastic gradient descent. The paper analyzes which data and model properties are responsible for the disproportionate impacts, why these aspects are affecting different groups disproportionately, and proposes guidelines to mitigate these effects. The proposed approach is evaluated on several datasets and settings.

[A Unified View of cGANs with and without Classifiers](#)

- Si-An Chen, Chun-Liang Li, Hsuan-Tien Lin
- abstract@[open-review](#): Conditional Generative Adversarial Networks (cGANs) are implicit generative models which allow to sample from class-conditional distributions. Existing cGANs are based on a wide range of different discriminator designs and training objectives. One popular design in earlier works is to include a classifier during training with the assumption that good classifiers can help eliminate samples generated with wrong classes. Nevertheless, including classifiers in cGANs often comes with a side effect of only generating easy-to-classify samples. Recently, some representative cGANs avoid the shortcoming and reach state-of-the-art performance without having classifiers. Somehow it remains unanswered whether the classifiers can be resurrected to design better cGANs. In this work, we demonstrate that classifiers can be properly leveraged to improve cGANs. We start by using the decomposition of the joint probability distribution to connect the goals of cGANs and classification as a unified framework. The framework, along with a classic energy model to parameterize distributions, justifies the use of classifiers for cGANs in a principled manner. It explains several popular cGAN variants, such as ACGAN, ProjGAN, and ContraGAN, as special cases with different levels of approximations, which provides a unified view and brings new insights to understanding cGANs. Experimental results demonstrate that the design inspired by the proposed framework outperforms state-of-the-art cGANs on multiple benchmark datasets, especially on the most challenging ImageNet. The code is available at <https://github.com/sian-chen/PyTorch-ECGAN>.

[Online and Offline Reinforcement Learning by Planning with a Learned Model](#)

- Julian Schrittwieser, Thomas Hubert, Amol Mandhane, Mohammadamin Barekatain, Ioannis Antonoglou, David Silver
- abstract@[open-review](#): Learning efficiently from small amounts of data has long been the focus of model-based reinforcement learning, both for the online case when interacting with the environment, and the offline case when learning from a fixed dataset. However, to date no single unified algorithm could demonstrate state-of-the-art results for both settings. In this work, we describe the Reanalyse algorithm, which uses model-based policy and value improvement operators to compute improved training targets for existing data points, allowing for efficient learning at data budgets varying by several orders of magnitude. We further show that Reanalyse can also be used to learn completely without environment interactions, as in the case of Offline Reinforcement Learning (Offline RL). Combining Reanalyse with the MuZero algorithm, we introduce MuZero Unplugged, a single unified algorithm for

any data budget, including Offline RL. In contrast to previous work, our algorithm requires no special adaptations for the off-policy or Offline RL settings. MuZero Unplugged sets new state-of-the-art results for Atari in the standard 200 million frame online setting as well as in the RL Unplugged Offline RL benchmark.

[Stochastic Multi-Armed Bandits with Control Variates](#)

- Arun Verma, Manjesh Kumar Hanawal
- abstract@[open-review](#): This paper studies a new variant of the stochastic multi-armed bandits problem where auxiliary information about the arm rewards is available in the form of control variates. In many applications like queuing and wireless networks, the arm rewards are functions of some exogenous variables. The mean values of these variables are known a priori from historical data and can be used as control variates. Leveraging the theory of control variates, we obtain mean estimates with smaller variance and tighter confidence bounds. We develop an upper confidence bound based algorithm named UCB-CV and characterize the regret bounds in terms of the correlation between rewards and control variates when they follow a multivariate normal distribution. We also extend UCB-CV to other distributions using resampling methods like Jackknifing and Splitting. Experiments on synthetic problem instances validate performance guarantees of the proposed algorithms.

[Near-Optimal No-Regret Learning in General Games](#)

- Constantinos Daskalakis, Maxwell Fishelson, Noah Golowich
- abstract@[open-review](#): We show that Optimistic Hedge -- a common variant of multiplicative-weights-updates with recency bias -- attains $\sqrt{\log T}$ regret in multi-player general-sum games. In particular, when every player of the game uses Optimistic Hedge to iteratively update her action in response to the history of play so far, then after T rounds of interaction, each player experiences total regret that is $\sqrt{\log T}$. Our bound improves, exponentially, the $O(T^{1/2})$ regret attainable by standard no-regret learners in games, the $O(T^{1/4})$ regret attainable by no-regret learners with recency bias (Syrgkanis et al., NeurIPS 2015), and the $O(T^{1/6})$ bound that was recently shown for Optimistic Hedge in the special case of two-player games (Chen & Peng, NeurIPS 2020). A direct corollary of our bound is that Optimistic Hedge converges to coarse correlated equilibrium in general games at a rate of $\tilde{O}(1/T)$.

[Improving Self-supervised Learning with Automated Unsupervised Outlier Arbitration](#)

- Yu Wang, Jingyang Lin, Jingjing Zou, Yingwei Pan, Ting Yao, Tao Mei
- abstract@[open-review](#): Our work reveals a structured shortcoming of the existing mainstream self-supervised learning methods. Whereas self-supervised learning frameworks usually take the prevailing perfect instance level invariance hypothesis for granted, we carefully investigate the pitfalls behind. Particularly, we argue that the existing augmentation pipeline for generating multiple positive views naturally introduces out-of-distribution (OOD) samples that undermine the learning of the downstream tasks. Generating diverse positive augmentations on the input does not always pay off in benefiting downstream tasks. To overcome this inherent deficiency, we introduce a lightweight latent variable model UOTA, targeting the view sampling issue for self-supervised learning. UOTA adaptively searches for the most important sampling region to produce views, and provides viable choice for outlier-robust self-supervised learning approaches. Our method directly generalizes to many mainstream self-supervised learning approaches, regardless of the loss's nature contrastive or not. We empirically show UOTA's advantage over the state-of-the-art self-supervised paradigms with evident margin, which well justifies the existence of the OOD sample issue embedded in the existing approaches. Especially, we theoretically prove that the merits of the proposal boil down to guaranteed estimator variance and bias reduction. Code is available: <https://github.com/ssl-codelab/uota>.

[Improving Anytime Prediction with Parallel Cascaded Networks and a Temporal-Difference Loss](#)

- Michael Iuzzolino, Michael C. Mozer, Samy Bengio
- abstract@[open-review](#): Although deep feedforward neural networks share some characteristics with the primate visual system, a key distinction is their dynamics. Deep nets typically operate in serial stages wherein each layer completes its computation before processing begins in subsequent layers. In contrast, biological systems have cascaded dynamics: information propagates from neurons at all layers in parallel but transmission occurs gradually over time, leading to speed-accuracy trade offs even in feedforward architectures. We explore the consequences of biologically inspired parallel hardware by constructing cascaded ResNets in which each residual block has propagation delays but all blocks update in parallel in a stateful manner. Because information transmitted through skip connections avoids delays, the functional depth of the architecture increases over time, yielding anytime predictions that improve with internal-processing time. We introduce a temporal-difference training loss that achieves a strictly superior speed-accuracy profile over standard losses and enables the cascaded architecture to outperform state-of-the-art anytime-prediction methods. The cascaded architecture has intriguing properties, including: it classifies typical instances more rapidly than atypical instances; it is more robust to both persistent and transient noise than is a conventional ResNet; and its time-varying output trace provides a signal that can be exploited to improve information processing and inference.

[Identifiable Generative models for Missing Not at Random Data Imputation](#)

- Chao Ma, Cheng Zhang
- abstract@[open-review](#): Real-world datasets often have missing values associated with complex generative processes, where the cause of the missingness may not be fully observed. This is known as missing not at random (MNAR) data. However, many imputation methods do not take into account the missingness mechanism, resulting in biased imputation values when MNAR data is present. Although there are a few methods that have considered the MNAR scenario, their model's identifiability under MNAR is generally not guaranteed. That is, model parameters can not be uniquely determined even with infinite data samples, hence the imputation results given by such models can still be biased. This issue is especially overlooked by many modern deep generative models. In this work, we fill in this gap by systematically analyzing the identifiability of generative models under MNAR. Furthermore, we propose a practical deep generative model which can provide identifiability guarantees under mild assumptions, for a wide range of MNAR mechanisms. Our method demonstrates a clear advantage for tasks on both synthetic data and multiple real-world scenarios with MNAR data.

[DNN-based Topology Optimisation: Spatial Invariance and Neural Tangent Kernel](#)

- Benjamin Dupuis, Arthur Jacot
- abstract@[open-review](#): We study the Solid Isotropic Material Penalization (SIMP) method with a density field generated by a fully-connected neural network, taking the coordinates as inputs. In the large width limit, we show that the use of DNNs leads to a filtering effect similar to traditional filtering techniques for SIMP, with a filter described by the Neural Tangent Kernel (NTK). This filter is however not invariant under translation, leading to visual artifacts and non-optimal shapes. We propose two embeddings of the input coordinates, which lead to (approximate) spatial invariance of the NTK and of the filter. We empirically confirm our theoretical observations and study how the filter size is affected by the architecture of the network. Our solution can easily be applied to any other coordinates-based generation method.

[Baleen: Robust Multi-Hop Reasoning at Scale via Condensed Retrieval](#)

- Omar Khattab, Christopher Potts, Matei Zaharia
- abstract@[open-review](#): Multi-hop reasoning (i.e., reasoning across two or more documents) is a key ingredient for NLP models that leverage large corpora to exhibit broad knowledge. To retrieve evidence passages, multi-hop models must contend with a fast-growing search space across the hops, represent

complex queries that combine multiple information needs, and resolve ambiguity about the best order in which to hop between training passages. We tackle these problems via Baleen, a system that improves the accuracy of multi-hop retrieval while learning robustly from weak training signals in the many-hop setting. To tame the search space, we propose condensed retrieval, a pipeline that summarizes the retrieved passages after each hop into a single compact context. To model complex queries, we introduce a focused late interaction retriever that allows different parts of the same query representation to match disparate relevant passages. Lastly, to infer the hopping dependencies among unordered training passages, we devise latent hop ordering, a weak-supervision strategy in which the trained retriever itself selects the sequence of hops. We evaluate Baleen on retrieval for two-hop question answering and many-hop claim verification, establishing state-of-the-art performance.

[Local Hyper-Flow Diffusion](#)

- Kimon Fountoulakis, Pan Li, Shenghao Yang
- abstract@[open-review](#): Recently, hypergraphs have attracted a lot of attention due to their ability to capture complex relations among entities. The resurgence of hypergraphs has resulted in data of increasing size and complexity that exhibit interesting small-scale and local structure, e.g., small-scale communities and localized node-ranking around a given set of seed nodes. Popular and principled ways to capture the local structure are the local hypergraph clustering problem and the related seed set expansion problem. In this work, we propose the first local diffusion method that achieves edge-size-independent Cheeger-type guarantee for the problem of local hypergraph clustering while applying to a rich class of higher-order relations that covers a number of previously studied special cases. Our method is based on a primal-dual optimization formulation where the primal problem has a natural network flow interpretation, and the dual problem has a cut-based interpretation using the ℓ_2 -norm penalty on associated cut-costs. We demonstrate the new technique is significantly better than state-of-the-art methods on both synthetic and real-world data.

[Permuton-induced Chinese Restaurant Process](#)

- Masahiro Nakano, Yasuhiro Fujiwara, Akisato Kimura, Takeshi Yamada, naonori ueda
- abstract@[open-review](#): This paper proposes the permuto-induced Chinese restaurant process (PCRP), a stochastic process on rectangular partitioning of a matrix. This distribution is suitable for use as a prior distribution in Bayesian nonparametric relational model to find hidden clusters in matrices and network data. Our main contribution is to introduce the notion of permutoins into the well-known Chinese restaurant process (CRP) for sequence partitioning: a permutoin is a probability measure on $[0,1]^{[0,1]}$ and can be regarded as a geometric interpretation of the scaling limit of permutations. Specifically, we extend the model that the table order of CRPs has a random geometric arrangement on $[0,1]^{[0,1]}$ drawn from the permutoin. By analogy with the relationship between the stick-breaking process (SBP) and CRP for the infinite mixture model of a sequence, this model can be regarded as a multi-dimensional extension of CRP paired with the block-breaking process (BBP), which has been recently proposed as a multi-dimensional extension of SBP. While BBP always has an infinite number of redundant intermediate variables, PCRP can be composed of varying size intermediate variables in a data-driven manner depending on the size and quality of the observation data. Experiments show that PCRP can improve the prediction performance in relational data analysis by reducing the local optima and slow mixing problems compared with the conventional BNP models because the local transitions of PCRP in Markov chain Monte Carlo inference are more flexible than the previous models.

[Faster Algorithms and Constant Lower Bounds for the Worst-Case Expected Error](#)

- Jonah Brown-Cohen
- abstract@[open-review](#): The study of statistical estimation without distributional assumptions on data values, but with knowledge of data collection methods was recently introduced by Chen, Valiant and Valiant (NeurIPS 2020). In this framework, the goal is to design estimators that minimize the worst-case expected error. Here the expectation is over a known, randomized data collection process from some population, and the data values corresponding to each element of the population are assumed to be worst-case. Chen, Valiant and Valiant show that, when data values are ℓ_∞ -normalized, there is a polynomial time algorithm to compute an estimator for the mean with worst-case expected error that is within a factor $\frac{1}{\sqrt{\pi}}$ of the optimum within the natural class of semilinear estimators. However, this algorithm is based on optimizing a somewhat complex concave objective function over a constrained set of positive semidefinite matrices, and thus does not come with explicit runtime guarantees beyond being polynomial time in the input. In this paper we design provably efficient algorithms for approximating the optimal semilinear estimator based on online convex optimization. In the setting where data values are ℓ_2 -normalized, our algorithm achieves a $\frac{1}{\sqrt{2}}$ -approximation by iteratively solving a sequence of standard SDPs. When data values are ℓ_2 -normalized, our algorithm iteratively computes the top eigenvector of a sequence of matrices, and does not lose any multiplicative approximation factor. Further, using experiments in settings where sample membership is correlated with data values (e.g. "importance sampling" and "snowball sampling"), we show that our ℓ_2 -normalized algorithm gives a similar advantage over standard estimators as the original ℓ_∞ -normalized algorithm of Chen, Valiant and Valiant, but with much lower computational complexity. We complement these positive results by stating a simple combinatorial condition which, if satisfied by a data collection process, implies that any (not necessarily semilinear) estimator for the mean has constant worst-case expected error.

[On Learning Domain-Invariant Representations for Transfer Learning with Multiple Sources](#)

- Trung Phung, Trung Le, Tung-Long Vuong, Toan Tran, Anh Tran, Hung Bui, Dinh Phung
- abstract@[open-review](#): Domain adaptation (DA) benefits from the rigorous theoretical works that study its insightful characteristics and various aspects, e.g., learning domain-invariant representations and its trade-off. However, it seems not the case for the multiple source DA and domain generalization (DG) settings which are remarkably more complicated and sophisticated due to the involvement of multiple source domains and potential unavailability of target domain during training. In this paper, we develop novel upper-bounds for the target general loss which appeal us to define two kinds of domain-invariant representations. We further study the pros and cons as well as the trade-offs of enforcing learning each domain-invariant representation. Finally, we conduct experiments to inspect the trade-off of these representations for offering practical hints regarding how to use them in practice and explore other interesting properties of our developed theory.

[You Never Cluster Alone](#)

- Yuming Shen, Ziyi Shen, Menghan Wang, Jie Qin, Philip Torr, Ling Shao
- abstract@[open-review](#): Recent advances in self-supervised learning with instance-level contrastive objectives facilitate unsupervised clustering. However, a standalone datum is not perceiving the context of the holistic cluster, and may undergo sub-optimal assignment. In this paper, we extend the mainstream contrastive learning paradigm to a cluster-level scheme, where all the data subjected to the same cluster contribute to a unified representation that encodes the context of each data group. Contrastive learning with this representation then rewards the assignment of each datum. To implement this vision, we propose twin-contrast clustering (TCC). We define a set of categorical variables as clustering assignment confidence, which links the instance-level learning track with the cluster-level one. On one hand, with the corresponding assignment variables being the weight, a weighted aggregation along the data points implements the set representation of a cluster. We further propose heuristic cluster augmentation equivalents to enable cluster-level contrastive learning. On the other hand, we derive the evidence lower-bound of the instance-level contrastive objective with the assignments. By reparametrizing the assignment variables, TCC is trained end-to-end, requiring no alternating steps. Extensive experiments show that TCC outperforms the state-of-the-art on benchmarked datasets.

[Dynamic COVID risk assessment accounting for community virus exposure from a spatial-temporal transmission model](#)

- Yuan Chen, Wenbo Fei, Qinxia Wang, Donglin Zeng, Yuanjia Wang

- abstract@[open-review](#): COVID-19 pandemic has caused unprecedented negative impacts on our society, including further exposing inequity and disparity in public health. To study the impact of socioeconomic factors on COVID transmission, we first propose a spatial-temporal model to examine the socioeconomic heterogeneity and spatial correlation of COVID-19 transmission at the community level. Second, to assess the individual risk of severe COVID-19 outcomes after a positive diagnosis, we propose a dynamic, varying-coefficient model that integrates individual-level risk factors from electronic health records (EHRs) with community-level risk factors. The underlying neighborhood prevalence of infections (both symptomatic and pre-symptomatic) predicted from the previous spatial-temporal model is included in the individual risk assessment so as to better capture the background risk of virus exposure for each individual. We design a weighting scheme to mitigate multiple selection biases inherited in EHRs of COVID patients. We analyze COVID transmission data in New York City (NYC, the epicenter of the first surge in the United States) and EHRs from NYC hospitals, where time-varying effects of community risk factors and significant interactions between individual- and community-level risk factors are detected. By examining the socioeconomic disparity of infection risks and interaction among the risk factors, our methods can assist public health decision-making and facilitate better clinical management of COVID patients.

[Dueling Bandits with Adversarial Sleeping](#)

- Aadirupa Saha, Pierre Gaillard
- abstract@[open-review](#): We introduce the problem of sleeping dueling bandits with stochastic preferences and adversarial availabilities (DB-SPAA). In almost all dueling bandit applications, the decision space often changes over time; eg, retail store management, online shopping, restaurant recommendation, search engine optimization, etc. Surprisingly, this 'sleeping aspect' of dueling bandits has never been studied in the literature. Like dueling bandits, the goal is to compete with the best arm by sequentially querying the preference feedback of item pairs. The non-triviality however results due to the non-stationary item spaces that allow any arbitrary subsets items to go unavailable every round. The goal is to find an optimal no-regret policy that can identify the best available item at each round, as opposed to the standard 'fixed best-arm regret objective' of dueling bandits. We first derive an instance-specific lower bound for DB-SPAA $\Omega(\sum_{i=1}^{K-1} \sum_{j=i+1}^K \frac{\log T}{\Delta(i,j)})$, where K is the number of items and $\Delta(i,j)$ is the gap between items i and j . This indicates that the sleeping problem with preference feedback is inherently more difficult than that for classical multi-armed bandits (MAB). We then propose two algorithms, with near optimal regret guarantees. Our results are corroborated empirically.

[Beware of the Simulated DAG! Causal Discovery Benchmarks May Be Easy to Game](#)

- Alexander Reisach, Christof Seiler, Sebastian Weichwald
- abstract@[open-review](#): Simulated DAG models may exhibit properties that, perhaps inadvertently, render their structure identifiable and unexpectedly affect structure learning algorithms. Here, we show that marginal variance tends to increase along the causal order for generically sampled additive noise models. We introduce varsorability as a measure of the agreement between the order of increasing marginal variance and the causal order. For commonly sampled graphs and model parameters, we show that the remarkable performance of some continuous structure learning algorithms can be explained by high varsorability and matched by a simple baseline method. Yet, this performance may not transfer to real-world data where varsorability may be moderate or dependent on the choice of measurement scales. On standardized data, the same algorithms fail to identify the ground-truth DAG or its Markov equivalence class. While standardization removes the pattern in marginal variance, we show that data generating processes that incur high varsorability also leave a distinct covariance pattern that may be exploited even after standardization. Our findings challenge the significance of generic benchmarks with independently drawn parameters. The code is available at <https://github.com/Scriddie/Varsortability>.

[Automated Dynamic Mechanism Design](#)

- Hanrui Zhang, Vincent Conitzer
- abstract@[open-review](#): We study Bayesian automated mechanism design in unstructured dynamic environments, where a principal repeatedly interacts with an agent, and takes actions based on the strategic agent's report of the current state of the world. Both the principal and the agent can have arbitrary and potentially different valuations for the actions taken, possibly also depending on the actual state of the world. Moreover, at any time, the state of the world may evolve arbitrarily depending on the action taken by the principal. The goal is to compute an optimal mechanism which maximizes the principal's utility in the face of the self-interested strategic agent. We give an efficient algorithm for computing optimal mechanisms, with or without payments, under different individual-rationality constraints, when the time horizon is constant. Our algorithm is based on a sophisticated linear program formulation, which can be customized in various ways to accommodate richer constraints. For environments with large time horizons, we show that the principal's optimal utility is hard to approximate within a certain constant factor, complementing our algorithmic result. These results paint a relatively complete picture for automated dynamic mechanism design in unstructured environments. We further consider a special case of the problem where the agent is myopic, and give a refined efficient algorithm whose time complexity scales linearly in the time horizon. In the full version of the paper, we show that memoryless mechanisms, which are without loss of generality optimal in Markov decision processes without strategic behavior, do not provide a good solution for our problem, in terms of both optimality and computational tractability. Moreover, we present experimental results where our algorithms are applied to synthetic dynamic environments with different characteristics, which not only serve as a proof of concept for our algorithms, but also exhibit intriguing phenomena in dynamic mechanism design.

[A generative nonparametric Bayesian model for whole genomes](#)

- Alan Amin, Eli N Weinstein, Debora Marks
- abstract@[open-review](#): Generative probabilistic modeling of biological sequences has widespread existing and potential use across biology and biomedicine, particularly given advances in high-throughput sequencing, synthesis and editing. However, we still lack methods with nucleotide resolution that are tractable at the scale of whole genomes and that can achieve high predictive accuracy in theory and practice. In this article we propose a new generative sequence model, the Bayesian embedded autoregressive (BEAR) model, which uses a parametric autoregressive model to specify a conjugate prior over a nonparametric Bayesian Markov model. We explore, theoretically and empirically, applications of BEAR models to a variety of statistical problems including density estimation, robust parameter estimation, goodness-of-fit tests, and two-sample tests. We prove rigorous asymptotic consistency results including nonparametric posterior concentration rates. We scale inference in BEAR models to datasets containing tens of billions of nucleotides. On genomic, transcriptomic, and metagenomic sequence data we show that BEAR models provide large increases in predictive performance as compared to parametric autoregressive models, among other results. BEAR models offer a flexible and scalable framework, with theoretical guarantees, for building and critiquing generative models at the whole genome scale.

[Robust Predictable Control](#)

- Ben Eysenbach, Russ R. Salakhutdinov, Sergey Levine
- abstract@[open-review](#): Many of the challenges facing today's reinforcement learning (RL) algorithms, such as robustness, generalization, transfer, and computational efficiency are closely related to compression. Prior work has convincingly argued why minimizing information is useful in the supervised learning setting, but standard RL algorithms lack an explicit mechanism for compression. The RL setting is unique because (1) its sequential nature allows an agent to use past information to avoid looking at future observations and (2) the agent can optimize its behavior to prefer states where decision making requires few bits. We take advantage of these properties to propose a method (RPC) for learning simple policies. This method brings together ideas from information bottlenecks, model-based RL, and bits-back coding into a simple and theoretically-justified algorithm. Our method jointly optimizes a latent-space model and policy to be self-consistent, such that the policy avoids states where the model is inaccurate. We demonstrate that our method achieves

much tighter compression than prior methods, achieving up to 5\$ \times \$ higher reward than a standard information bottleneck when constrained to use just 0.3 bits per observation. We also demonstrate that our method learns policies that are more robust and generalize better to new tasks.

[Unsupervised Speech Recognition](#)

- Alexei Baevski, Wei-Ning Hsu, Alexis CONNEAU, Michael Auli
- abstract@[open-review](#): Despite rapid progress in the recent past, current speech recognition systems still require labeled training data which limits this technology to a small fraction of the languages spoken around the globe. This paper describes wav2vec-U, short for wav2vec Unsupervised, a method to train speech recognition models without any labeled data. We leverage self-supervised speech representations to segment unlabeled audio and learn a mapping from these representations to phonemes via adversarial training. The right representations are key to the success of our method. Compared to the best previous unsupervised work, wav2vec-U reduces the phone error rate on the TIMIT benchmark from 26.1 to 11.3. On the larger English LibriSpeech benchmark, wav2vec-U achieves a word error rate of 5.9 on test-other, rivaling some of the best published systems trained on 960 hours of labeled data from only two years ago. We also experiment on nine other languages, including low-resource languages such as Kyrgyz, Swahili and Tatar.

[Robustness between the worst and average case](#)

- Leslie Rice, Anna Bair, Huan Zhang, J. Zico Kolter
- abstract@[open-review](#): Several recent works in machine learning have focused on evaluating the test-time robustness of a classifier: how well the classifier performs not just on the target domain it was trained upon, but upon perturbed examples. In these settings, the focus has largely been on two extremes of robustness: the robustness to perturbations drawn at random from within some distribution (i.e., robustness to random perturbations), and the robustness to the worst case perturbation in some set (i.e., adversarial robustness). In this paper, we argue that a sliding scale between these two extremes provides a valuable additional metric by which to gauge robustness. Specifically, we illustrate that each of these two extremes is naturally characterized by a (functional) q -norm over perturbation space, with $q=1$ corresponding to robustness to random perturbations and $q=\infty$ corresponding to adversarial perturbations. We then present the main technical contribution of our paper: a method for efficiently estimating the value of these norms by interpreting them as the partition function of a particular distribution, then using path sampling with MCMC methods to estimate this partition function (either traditional Metropolis-Hastings for non-differentiable perturbations, or Hamiltonian Monte Carlo for differentiable perturbations). We show that our approach provides substantially better estimates than simple random sampling of the actual “intermediate- q ” robustness of both standard, data-augmented, and adversarially-trained classifiers, illustrating a clear tradeoff between classifiers that optimize different metrics. Code for reproducing experiments can be found at https://github.com/locuslab/intermediate_robustness.

[Online Learning and Control of Complex Dynamical Systems from Sensory Input](#)

- Oumayma Bounou, Jean Ponce, Justin Carpenter
- abstract@[open-review](#): Identifying an effective model of a dynamical system from sensory data and using it for future state prediction and control is challenging. Recent data-driven algorithms based on Koopman theory are a promising approach to this problem, but they typically never update the model once it has been identified from a relatively small set of observations, thus making long-term prediction and control difficult for realistic systems, in robotics or fluid mechanics for example. This paper introduces a novel method for learning an embedding of the state space with linear dynamics from sensory data. Unlike previous approaches, the dynamics model can be updated online and thus easily applied to systems with non-linear dynamics in the original configuration space. The proposed approach is evaluated empirically on several classical dynamical systems and sensory modalities, with good performance on long-term prediction and control.

[Self-Supervised Bug Detection and Repair](#)

- Miltiadis Allamanis, Henry Jackson-Flux, Marc Brockschmidt
- abstract@[open-review](#): Machine learning-based program analyses have recently shown the promise of integrating formal and probabilistic reasoning towards aiding software development. However, in the absence of large annotated corpora, training these analyses is challenging. Towards addressing this, we present BugLab, an approach for self-supervised learning of bug detection and repair. BugLab co-trains two models: (1) a detector model that learns to detect and repair bugs in code, (2) a selector model that learns to create buggy code for the detector to use as training data. A Python implementation of BugLab improves by 30% upon baseline methods on a test dataset of 2374 real-life bugs and finds 19 previously unknown bugs in open-source software.

[Faster Neural Network Training with Approximate Tensor Operations](#)

- Menachem Adelman, Kfir Levy, Ido Hakimi, Mark Silberstein
- abstract@[open-review](#): We propose a novel technique for faster deep neural network training which systematically applies sample-based approximation to the constituent tensor operations, i.e., matrix multiplications and convolutions. We introduce new sampling techniques, study their theoretical properties, and prove that they provide the same convergence guarantees when applied to SGD training. We apply approximate tensor operations to single and multi-node training of MLP and CNN networks on MNIST, CIFAR-10 and ImageNet datasets. We demonstrate up to 66% reduction in the amount of computations and communication, and up to 1.37x faster training time while maintaining negligible or no impact on the final test accuracy.

[Learning Interpretable Decision Rule Sets: A Submodular Optimization Approach](#)

- Fan Yang, Kai He, Linxiao Yang, Hongxia Du, Jingbang Yang, Bo Yang, Liang Sun
- abstract@[open-review](#): Rule sets are highly interpretable logical models in which the predicates for decision are expressed in disjunctive normal form (DNF, OR-of-ANDs), or, equivalently, the overall model comprises an unordered collection of if-then decision rules. In this paper, we consider a submodular optimization based approach for learning rule sets. The learning problem is framed as a subset selection task in which a subset of all possible rules needs to be selected to form an accurate and interpretable rule set. We employ an objective function that exhibits submodularity and thus is amenable to submodular optimization techniques. To overcome the difficulty arose from dealing with the exponential-sized ground set of rules, the subproblem of searching a rule is casted as another subset selection task that asks for a subset of features. We show it is possible to write the induced objective function for the subproblem as a difference of two submodular (DS) functions to make it approximately solvable by DS optimization algorithms. Overall, the proposed approach is simple, scalable, and likely to be benefited from further research on submodular optimization. Experiments on real datasets demonstrate the effectiveness of our method.

[Spatial-Temporal Super-Resolution of Satellite Imagery via Conditional Pixel Synthesis](#)

- Yutong He, Dingjie Wang, Nicholas Lai, William Zhang, Chenlin Meng, Marshall Burke, David Lobell, Stefano Ermon
- abstract@[open-review](#): High-resolution satellite imagery has proven useful for a broad range of tasks, including measurement of global human population, local economic livelihoods, and biodiversity, among many others. Unfortunately, high-resolution imagery is both infrequently collected and expensive to purchase, making it hard to efficiently and effectively scale these downstream tasks over both time and space. We propose a new conditional pixel synthesis model that uses abundant, low-cost, low-resolution imagery to generate accurate high-resolution imagery at locations and times in which it is unavailable. We show that our model attains photo-realistic sample quality and outperforms competing baselines on a key downstream task – object counting – particularly in geographic locations where conditions on the ground are changing rapidly.

On Memorization in Probabilistic Deep Generative Models

- Gerrit van den Burg, Chris Williams
- abstract@[open-review](#): Recent advances in deep generative models have led to impressive results in a variety of application domains. Motivated by the possibility that deep learning models might memorize part of the input data, there have been increased efforts to understand how memorization arises. In this work, we extend a recently proposed measure of memorization for supervised learning (Feldman, 2019) to the unsupervised density estimation problem and adapt it to be more computationally efficient. Next, we present a study that demonstrates how memorization can occur in probabilistic deep generative models such as variational autoencoders. This reveals that the form of memorization to which these models are susceptible differs fundamentally from mode collapse and overfitting. Furthermore, we show that the proposed memorization score measures a phenomenon that is not captured by commonly-used nearest neighbor tests. Finally, we discuss several strategies that can be used to limit memorization in practice. Our work thus provides a framework for understanding problematic memorization in probabilistic generative models.

You Are the Best Reviewer of Your Own Papers: An Owner-Assisted Scoring Mechanism

- Weijie Su
- abstract@[open-review](#): I consider the setting where reviewers offer very noisy scores for a number of items for the selection of high-quality ones (e.g., peer review of large conference proceedings) whereas the owner of these items knows the true underlying scores but prefers not to provide this information. To address this withholding of information, in this paper, I introduce the Isotonic Mechanism, a simple and efficient approach to improving on the imprecise raw scores by leveraging certain information that the owner is incentivized to provide. This mechanism takes as input the ranking of the items from best to worst provided by the owner, in addition to the raw scores provided by the reviewers. It reports adjusted scores for the items by solving a convex optimization problem. Under certain conditions, I show that the owner's optimal strategy is to honestly report the true ranking of the items to her best knowledge in order to maximize the expected utility. Moreover, I prove that the adjusted scores provided by this owner-assisted mechanism are indeed significantly more accurate than the raw scores provided by the reviewers. This paper concludes with several extensions of the Isotonic Mechanism and some refinements of the mechanism for practical considerations.

Garment4D: Garment Reconstruction from Point Cloud Sequences

- Fangzhou Hong, Liang Pan, Zhongang Cai, Ziwei Liu
- abstract@[open-review](#): Learning to reconstruct 3D garments is important for dressing 3D human bodies of different shapes in different poses. Previous works typically rely on 2D images as input, which however suffer from the scale and pose ambiguities. To circumvent the problems caused by 2D images, we propose a principled framework, Garment4D, that uses 3D point cloud sequences of dressed humans for garment reconstruction. Garment4D has three dedicated steps: sequential garments registration, canonical garment estimation, and posed garment reconstruction. The main challenges are two-fold: 1) effective 3D feature learning for fine details, and 2) capture of garment dynamics caused by the interaction between garments and the human body, especially for loose garments like skirts. To unravel these problems, we introduce a novel Proposal-Guided Hierarchical Feature Network and Iterative Graph Convolution Network, which integrate both high-level semantic features and low-level geometric features for fine details reconstruction. Furthermore, we propose a Temporal Transformer for smooth garment motions capture. Unlike non-parametric methods, the reconstructed garment meshes by our method are separable from the human body and have strong interpretability, which is desirable for downstream tasks. As the first attempt at this task, high-quality reconstruction results are qualitatively and quantitatively illustrated through extensive experiments. Codes are available at <https://github.com/hongfz16/Garment4D>.

Fast Policy Extragradient Methods for Competitive Games with Entropy Regularization

- Shicong Cen, Yuting Wei, Yuejie Chi
- abstract@[open-review](#): This paper investigates the problem of computing the equilibrium of competitive games, which is often modeled as a constrained saddle-point optimization problem with probability simplex constraints. Despite recent efforts in understanding the last-iterate convergence of extragradient methods in the unconstrained setting, the theoretical underpinnings of these methods in the constrained settings, especially those using multiplicative updates, remain highly inadequate, even when the objective function is bilinear. Motivated by the algorithmic role of entropy regularization in single-agent reinforcement learning and game theory, we develop provably efficient extragradient methods to find the quantal response equilibrium (QRE)---which are solutions to zero-sum two-player matrix games with entropy regularization---at a linear rate. The proposed algorithms can be implemented in a decentralized manner, where each player executes symmetric and multiplicative updates iteratively using its own payoff without observing the opponent's actions directly. In addition, by controlling the knob of entropy regularization, the proposed algorithms can locate an approximate Nash equilibrium of the unregularized matrix game at a sublinear rate without assuming the Nash equilibrium to be unique. Our methods also lead to efficient policy extragradient algorithms for solving entropy-regularized zero-sum Markov games at a linear rate. All of our convergence rates are nearly dimension-free, which are independent of the size of the state and action spaces up to logarithm factors, highlighting the positive role of entropy regularization for accelerating convergence.

Shift-Robust GNNs: Overcoming the Limitations of Localized Graph Training data

- Qi Zhu, Natalia Ponomareva, Jiawei Han, Bryan Perozzi
- abstract@[open-review](#): There has been a recent surge of interest in designing Graph Neural Networks (GNNs) for semi-supervised learning tasks. Unfortunately this work has assumed that the nodes labeled for use in training were selected uniformly at random (i.e. are an IID sample). However in many real world scenarios gathering labels for graph nodes is both expensive and inherently biased -- so this assumption can not be met. GNNs can suffer poor generalization when this occurs, by overfitting to superfluous regularities present in the training data. In this work we present a method, Shift-Robust GNN (SR-GNN), designed to account for distributional differences between biased training data and the graph's true inference distribution. SR-GNN adapts GNN models for the presence of distributional shifts between the nodes which have had labels provided for training and the rest of the dataset. We illustrate the effectiveness of SR-GNN in a variety of experiments with biased training datasets on common GNN benchmark datasets for semi-supervised learning, where we see that SR-GNN outperforms other GNN baselines by accuracy, eliminating at least (~40%) of the negative effects introduced by biased training data. On the largest dataset we consider, ogb-arxiv, we observe an 2% absolute improvement over the baseline and reduce 30% of the negative effects.

RIM: Reliable Influence-based Active Learning on Graphs

- Wentao Zhang, Yixin Wang, Zhenbang You, Meng Cao, Ping Huang, Jilong Shan, Zhi Yang, Bin CUI
- abstract@[open-review](#): Message passing is the core of most graph models such as Graph Convolutional Network (GCN) and Label Propagation (LP), which usually require a large number of clean labeled data to smooth out the neighborhood over the graph. However, the labeling process can be tedious, costly, and error-prone in practice. In this paper, we propose to unify active learning (AL) and message passing towards minimizing labeling costs, e.g., making use of few and unreliable labels that can be obtained cheaply. We make two contributions towards that end. First, we open up a perspective by drawing a connection between AL enforcing message passing and social influence maximization, ensuring that the selected samples effectively improve the model performance. Second, we propose an extension to the influence model that incorporates an explicit quality factor to model label noise. In this way, we derive a fundamentally new AL selection criterion for GCN and LP--reliable influence maximization (RIM)--by considering quantity and quality of influence simultaneously. Empirical studies on public datasets show that RIM significantly outperforms current AL methods in terms of accuracy and efficiency.

[Dynamical Wasserstein Barycenters for Time-series Modeling](#)

- Kevin Cheng, Shuchin Aeron, Michael C. Hughes, Eric L Miller
- abstract@[open-review](#): Many time series can be modeled as a sequence of segments representing high-level discrete states, such as running and walking in a human activity application. Flexible models should describe the system state and observations in stationary ``pure-state'' periods as well as transition periods between adjacent segments, such as a gradual slowdown between running and walking. However, most prior work assumes instantaneous transitions between pure discrete states. We propose a dynamical Wasserstein barycentric (DWB) model that estimates the system state over time as well as the data-generating distributions of pure states in an unsupervised manner. Our model assumes each pure state generates data from a multivariate normal distribution, and characterizes transitions between states via displacement-interpolation specified by the Wasserstein barycenter. The system state is represented by a barycentric weight vector which evolves over time via a random walk on the simplex. Parameter learning leverages the natural Riemannian geometry of Gaussian distributions under the Wasserstein distance, which leads to improved convergence speeds. Experiments on several human activity datasets show that our proposed DWB model accurately learns the generating distribution of pure states while improving state estimation for transition periods compared to the commonly used linear interpolation mixture models.

[RelaySum for Decentralized Deep Learning on Heterogeneous Data](#)

- Thijs Vogels, Lie He, Anastasiia Koloskova, Sai Praneeth Karimireddy, Tao Lin, Sebastian U. Stich, Martin Jaggi
- abstract@[open-review](#): In decentralized machine learning, workers compute model updates on their local data. Because the workers only communicate with few neighbors without central coordination, these updates propagate progressively over the network. This paradigm enables distributed training on networks without all-to-all connectivity, helping to protect data privacy as well as to reduce the communication cost of distributed training in data centers. A key challenge, primarily in decentralized deep learning, remains the handling of differences between the workers' local data distributions. To tackle this challenge, we introduce the RelaySum mechanism for information propagation in decentralized learning. RelaySum uses spanning trees to distribute information exactly uniformly across all workers with finite delays depending on the distance between nodes. In contrast, the typical gossip averaging mechanism only distributes data uniformly asymptotically while using the same communication volume per step as RelaySum. We prove that RelaySGD, based on this mechanism, is independent of data heterogeneity and scales to many workers, enabling highly accurate decentralized deep learning on heterogeneous data.

[Transformers Generalize DeepSets and Can be Extended to Graphs & Hypergraphs](#)

- Jinwoo Kim, Saeyoon Oh, Seunghoon Hong
- abstract@[open-review](#): We present a generalization of Transformers to any-order permutation invariant data (sets, graphs, and hypergraphs). We begin by observing that Transformers generalize DeepSets, or first-order (set-input) permutation invariant MLPs. Then, based on recently characterized higher-order invariant MLPs, we extend the concept of self-attention to higher orders and propose higher-order Transformers for order-\$k\$ data (\$k=2\$ for graphs and \$k>2\$ for hypergraphs). Unfortunately, higher-order Transformers turn out to have prohibitive complexity \$\mathcal{O}(n^{2k})\$ to the number of input nodes \$n\$. To address this problem, we present sparse higher-order Transformers that have quadratic complexity to the number of input hyperedges, and further adopt the kernel attention approach to reduce the complexity to linear. In particular, we show that the sparse second-order Transformers with kernel attention are theoretically more expressive than message passing operations while having an asymptotically identical complexity. Our models achieve significant performance improvement over invariant MLPs and message-passing graph neural networks in large-scale graph regression and set-to-(hyper)graph prediction tasks. Our implementation is available at <https://github.com/jw9730/hot>.

[No Regrets for Learning the Prior in Bandits](#)

- Soumya Basu, Branislav Kveton, Manzil Zaheer, Csaba Szepesvari
- abstract@[open-review](#): We propose AdaTS, a Thompson sampling algorithm that adapts sequentially to bandit tasks that it interacts with. The key idea in AdaTS is to adapt to an unknown task prior distribution by maintaining a distribution over its parameters. When solving a bandit task, that uncertainty is marginalized out and properly accounted for. AdaTS is a fully-Bayesian algorithm that can be implemented efficiently in several classes of bandit problems. We derive upper bounds on its Bayes regret that quantify the loss due to not knowing the task prior, and show that it is small. Our theory is supported by experiments, where AdaTS outperforms prior algorithms and works well even in challenging real-world problems.

[Encoding Robustness to Image Style via Adversarial Feature Perturbations](#)

- Manli Shu, Zuxuan Wu, Micah Goldblum, Tom Goldstein
- abstract@[open-review](#): Adversarial training is the industry standard for producing models that are robust to small adversarial perturbations. However, machine learning practitioners need models that are robust to other kinds of changes that occur naturally, such as changes in the style or illumination of input images. Such changes in input distribution have been effectively modeled as shifts in the mean and variance of deep image features. We adapt adversarial training by directly perturbing feature statistics, rather than image pixels, to produce models that are robust to various unseen distributional shifts. We explore the relationship between these perturbations and distributional shifts by visualizing adversarial features. Our proposed method, Adversarial Batch Normalization (AdvBN), is a single network layer that generates worst-case feature perturbations during training. By fine-tuning neural networks on adversarial feature distributions, we observe improved robustness of networks to various unseen distributional shifts, including style variations and image corruptions. In addition, we show that our proposed adversarial feature perturbation can be complementary to existing image space data augmentation methods, leading to improved performance. The source code and pre-trained models are released at <https://github.com/azshue/AdvBN>.

[Continuized Accelerations of Deterministic and Stochastic Gradient Descents, and of Gossip Algorithms](#)

- Mathieu Even, Raphaël Berthier, Francis Bach, Nicolas Flammarion, Hadrien Hendrikx, Pierre Gaillard, Laurent Massoulié, Adrien Taylor
- abstract@[open-review](#): We introduce the ``continuized'' Nesterov acceleration, a close variant of Nesterov acceleration whose variables are indexed by a continuous time parameter. The two variables continuously mix following a linear ordinary differential equation and take gradient steps at random times. This continuized variant benefits from the best of the continuous and the discrete frameworks: as a continuous process, one can use differential calculus to analyze convergence and obtain analytical expressions for the parameters; but a discretization of the continuized process can be computed exactly with convergence rates similar to those of Nesterov original acceleration. We show that the discretization has the same structure as Nesterov acceleration, but with random parameters. We provide continuized Nesterov acceleration under deterministic as well as stochastic gradients, with either additive or multiplicative noise. Finally, using our continuized framework and expressing the gossip averaging problem as the stochastic minimization of a certain energy function, we provide the first rigorous acceleration of asynchronous gossip algorithms.

[Natural continual learning: success is a journey, not \(just\) a destination](#)

- Ta-Chu Kao, Kristopher Jensen, Gido van de Ven, Alberto Bernacchia, Guillaume Hennequin
- abstract@[open-review](#): Biological agents are known to learn many different tasks over the course of their lives, and to be able to revisit previous tasks and behaviors with little to no loss in performance. In contrast, artificial agents are prone to ‘catastrophic forgetting’ whereby performance on previous tasks deteriorates rapidly as new ones are acquired. This shortcoming has recently been addressed using methods that encourage parameters to stay close to

those used for previous tasks. This can be done by (i) using specific parameter regularizers that map out suitable destinations in parameter space, or (ii) guiding the optimization journey by projecting gradients into subspaces that do not interfere with previous tasks. However, these methods often exhibit subpar performance in both feedforward and recurrent neural networks, with recurrent networks being of interest to the study of neural dynamics supporting biological continual learning. In this work, we propose Natural Continual Learning (NCL), a new method that unifies weight regularization and projected gradient descent. NCL uses Bayesian weight regularization to encourage good performance on all tasks at convergence and combines this with gradient projection using the prior precision, which prevents catastrophic forgetting during optimization. Our method outperforms both standard weight regularization techniques and projection based approaches when applied to continual learning problems in feedforward and recurrent networks. Finally, the trained networks evolve task-specific dynamics that are strongly preserved as new tasks are learned, similar to experimental findings in biological circuits.

[Individual Privacy Accounting via a Rényi Filter](#)

- Vitaly Feldman, Tijana Zrnic
- abstract@[open-review](#): We consider a sequential setting in which a single dataset of individuals is used to perform adaptively-chosen analyses, while ensuring that the differential privacy loss of each participant does not exceed a pre-specified privacy budget. The standard approach to this problem relies on bounding a worst-case estimate of the privacy loss over all individuals and all possible values of their data, for every single analysis. Yet, in many scenarios this approach is overly conservative, especially for "typical" data points which incur little privacy loss by participation in most of the analyses. In this work, we give a method for tighter privacy loss accounting based on the value of a personalized privacy loss estimate for each individual in each analysis. To implement the accounting method we design a filter for Rényi differential privacy. A filter is a tool that ensures that the privacy parameter of a composed sequence of algorithms with adaptively-chosen privacy parameters does not exceed a pre-specified budget. Our filter is simpler and tighter than the known filter for (ϵ, δ) -differential privacy by Rogers et al. (2016). We apply our results to the analysis of noisy gradient descent and show that personalized accounting can be practical, easy to implement, and can only make the privacy-utility tradeoff tighter.

[Post-Training Quantization for Vision Transformer](#)

- Zhenhua Liu, Yunhe Wang, Kai Han, Wei Zhang, Siwei Ma, Wen Gao
- abstract@[open-review](#): Recently, transformer has achieved remarkable performance on a variety of computer vision applications. Compared with mainstream convolutional neural networks, vision transformers are often of sophisticated architectures for extracting powerful feature representations, which are more difficult to be developed on mobile devices. In this paper, we present an effective post-training quantization algorithm for reducing the memory storage and computational costs of vision transformers. Basically, the quantization task can be regarded as finding the optimal low-bit quantization intervals for weights and inputs, respectively. To preserve the functionality of the attention mechanism, we introduce a ranking loss into the conventional quantization objective that aims to keep the relative order of the self-attention results after quantization. Moreover, we thoroughly analyze the relationship between quantization loss of different layers and the feature diversity, and explore a mixed-precision quantization scheme by exploiting the nuclear norm of each attention map and output feature. The effectiveness of the proposed method is verified on several benchmark models and datasets, which outperforms the state-of-the-art post-training quantization algorithms. For instance, we can obtain an 81.29% top-1 accuracy using DeiT-B model on ImageNet dataset with about 8-bit quantization. Code will be available at <https://gitee.com/mindspore/models/tree/master/research/cv/VT-PTQ>.

[Unsupervised Part Discovery from Contrastive Reconstruction](#)

- Subhabrata Choudhury, Iro Laina, Christian Rupprecht, Andrea Vedaldi
- abstract@[open-review](#): The goal of self-supervised visual representation learning is to learn strong, transferable image representations, with the majority of research focusing on object or scene level. On the other hand, representation learning at part level has received significantly less attention. In this paper, we propose an unsupervised approach to object part discovery and segmentation and make three contributions. First, we construct a proxy task through a set of objectives that encourages the model to learn a meaningful decomposition of the image into its parts. Secondly, prior work argues for reconstructing or clustering pre-computed features as a proxy to parts; we show empirically that this alone is unlikely to find meaningful parts; mainly because of their low resolution and the tendency of classification networks to spatially smear out information. We suggest that image reconstruction at the level of pixels can alleviate this problem, acting as a complementary cue. Lastly, we show that the standard evaluation based on keypoint regression does not correlate well with segmentation quality and thus introduce different metrics, NMI and ARI, that better characterize the decomposition of objects into parts. Our method yields semantic parts which are consistent across fine-grained but visually distinct categories, outperforming the state of the art on three benchmark datasets. Code is available at the project page: <https://www.robots.ox.ac.uk/~vgg/research/unsup-parts/>.

[ASSANet: An Anisotropic Separable Set Abstraction for Efficient Point Cloud Representation Learning](#)

- Guocheng Qian, Hasan Hammoud, Guohao Li, Ali Thabet, Bernard Ghanem
- abstract@[open-review](#): Access to 3D point cloud representations has been widely facilitated by LiDAR sensors embedded in various mobile devices. This has led to an emerging need for fast and accurate point cloud processing techniques. In this paper, we revisit and dive deeper into PointNet++, one of the most influential yet under-explored networks, and develop faster and more accurate variants of the model. We first present a novel Separable Set Abstraction (SA) module that disentangles the vanilla SA module used in PointNet++ into two separate learning stages: (1) learning channel correlation and (2) learning spatial correlation. The Separable SA module is significantly faster than the vanilla version, yet it achieves comparable performance. We then introduce a new Anisotropic Reduction function into our Separable SA module and propose an Anisotropic Separable SA (ASSA) module that substantially increases the network's accuracy. We later replace the vanilla SA modules in PointNet++ with the proposed ASSA modules, and denote the modified network as ASSANet. Extensive experiments on point cloud classification, semantic segmentation, and part segmentation show that ASSANet outperforms PointNet++ and other methods, achieving much higher accuracy and faster speeds. In particular, ASSANet outperforms PointNet++ by 7.4% mIoU on S3DIS Area 5, while maintaining 1.6 times faster inference speed on a single NVIDIA 2080Ti GPU. Our scaled ASSANet variant achieves 66.8% mIoU and outperforms KPConv, while being more than 54 times faster.

[An Empirical Investigation of Domain Generalization with Empirical Risk Minimizers](#)

- Ramakrishna Vedantam, David Lopez-Paz, David J. Schwab
- abstract@[open-review](#): Recent work demonstrates that deep neural networks trained using Empirical Risk Minimization (ERM) can generalize under distribution shift, outperforming specialized training algorithms for domain generalization. The goal of this paper is to further understand this phenomenon. In particular, we study the extent to which the seminal domain adaptation theory of Ben-David et al. (2007) explains the performance of ERMs. Perhaps surprisingly, we find that this theory does not provide a tight explanation of the out-of-domain generalization observed across a large number of ERM models trained on three popular domain generalization datasets. This motivates us to investigate other possible measures—that, however, lack theory—which could explain generalization in this setting. Our investigation reveals that measures relating to the Fisher information, predictive entropy, and maximum mean discrepancy are good predictors of the out-of-distribution generalization of ERM models. We hope that our work helps galvanize the community towards building a better understanding of when deep networks trained with ERM generalize out-of-distribution.

[Fair Sequential Selection Using Supervised Learning Models](#)

- Mohammad Mahdi Khalili, Xueru Zhang, Mahed Abroshan

- abstract@[open-review](#): We consider a selection problem where sequentially arrived applicants apply for a limited number of positions/jobs. At each time step, a decision maker accepts or rejects the given applicant using a pre-trained supervised learning model until all the vacant positions are filled. In this paper, we discuss whether the fairness notions (e.g., equal opportunity, statistical parity, etc.) that are commonly used in classification problems are suitable for the sequential selection problems. In particular, we show that even with a pre-trained model that satisfies the common fairness notions, the selection outcomes may still be biased against certain demographic groups. This observation implies that the fairness notions used in classification problems are not suitable for a selection problem where the applicants compete for a limited number of positions. We introduce a new fairness notion, "Equal Selection (ES)," suitable for sequential selection problems and propose a post-processing approach to satisfy the ES fairness notion. We also consider a setting where the applicants have privacy concerns, and the decision maker only has access to the noisy version of sensitive attributes. In this setting, we can show that the \textit{perfect} ES fairness can still be attained under certain conditions.

[Towards Sample-efficient Overparameterized Meta-learning](#)

- Yue Sun, Adhyyan Narang, Ibrahim Gulluk, Samet Oymak, Maryam Fazel
- abstract@[open-review](#): An overarching goal in machine learning is to build a generalizable model with few samples. To this end, overparameterization has been the subject of immense interest to explain the generalization ability of deep nets even when the size of the dataset is smaller than that of the model. While the prior literature focuses on the classical supervised setting, this paper aims to demystify overparameterization for meta-learning. Here we have a sequence of linear-regression tasks and we ask: (1) Given earlier tasks, what is the optimal linear representation of features for a new downstream task? and (2) How many samples do we need to build this representation? This work shows that surprisingly, overparameterization arises as a natural answer to these fundamental meta-learning questions. Specifically, for (1), we first show that learning the optimal representation coincides with the problem of designing a task-aware regularization to promote inductive bias. We leverage this inductive bias to explain how the downstream task actually benefits from overparameterization, in contrast to prior works on few-shot learning. For (2), we develop a theory to explain how feature covariance can implicitly help reduce the sample complexity well below the degrees of freedom and lead to small estimation error. We then integrate these findings to obtain an overall performance guarantee for our meta-learning algorithm. Numerical experiments on real and synthetic data verify our insights on overparameterized meta-learning.

[ScaleCert: Scalable Certified Defense against Adversarial Patches with Sparse Superficial Layers](#)

- Husheng Han, Kaidi Xu, Xing Hu, Xiaobing Chen, Ling LIANG, Zidong Du, Qi Guo, Yanzhi Wang, Yunji Chen
- abstract@[open-review](#): Adversarial patch attacks that craft the pixels in a confined region of the input images show their powerful attack effectiveness in physical environments even with noises or deformations. Existing certified defenses towards adversarial patch attacks work well on small images like MNIST and CIFAR-10 datasets, but achieve very poor certified accuracy on higher-resolution images like ImageNet. It is urgent to design both robust and effective defenses against such a practical and harmful attack in industry-level larger images. In this work, we propose the certified defense methodology that achieves high provable robustness for high-resolution images and largely improves the practicality for real adoption of the certified defense. The basic insight of our work is that the adversarial patch intends to leverage localized superficial important neurons (SIN) to manipulate the prediction results. Hence, we leverage the SIN-based DNN compression techniques to significantly improve the certified accuracy, by reducing the adversarial region searching overhead and filtering the prediction noises. Our experimental results show that the certified accuracy is increased from 36.3% (the state-of-the-art certified detection) to 60.4% on the ImageNet dataset, largely pushing the certified defenses for practical use.

[Towards mental time travel: a hierarchical memory for reinforcement learning agents](#)

- Andrew Lampinen, Stephanie Chan, Andrea Banino, Felix Hill
- abstract@[open-review](#): Reinforcement learning agents often forget details of the past, especially after delays or distractor tasks. Agents with common memory architectures struggle to recall and integrate across multiple timesteps of a past event, or even to recall the details of a single timestep that is followed by distractor tasks. To address these limitations, we propose a Hierarchical Chunk Attention Memory (HCAM), that helps agents to remember the past in detail. HCAM stores memories by dividing the past into chunks, and recalls by first performing high-level attention over coarse summaries of the chunks, and then performing detailed attention within only the most relevant chunks. An agent with HCAM can therefore "mentally time-travel"--remember past events in detail without attending to all intervening events. We show that agents with HCAM substantially outperform agents with other memory architectures at tasks requiring long-term recall, retention, or reasoning over memory. These include recalling where an object is hidden in a 3D environment, rapidly learning to navigate efficiently in a new neighborhood, and rapidly learning and retaining new words. Agents with HCAM can extrapolate to task sequences much longer than they were trained on, and can even generalize zero-shot from a meta-learning setting to maintaining knowledge across episodes. HCAM improve agent sample efficiency, generalization, and generality (by solving tasks that previously required specialized architectures). Our work is a step towards agents that can learn, interact, and adapt in complex and temporally-extended environments.

[Beyond Tikhonov: faster learning with self-concordant losses, via iterative regularization](#)

- Gaspard Beugnot, Julien Mairal, Alessandro Rudi
- abstract@[open-review](#): The theory of spectral filtering is a remarkable tool to understand the statistical properties of learning with kernels. For least squares, it allows to derive various regularization schemes that yield faster convergence rates of the excess risk than with Tikhonov regularization. This is typically achieved by leveraging classical assumptions called source and capacity conditions, which characterize the difficulty of the learning task. In order to understand estimators derived from other loss functions, Marteau-Ferey et al. have extended the theory of Tikhonov regularization to generalized self concordant loss functions (GSC), which contain, e.g., the logistic loss. In this paper, we go a step further and show that fast and optimal rates can be achieved for GSC by using the iterated Tikhonov regularization scheme, which is intrinsically related to the proximal point method in optimization, and overcomes the limitation of the classical Tikhonov regularization.

[Variational Bayesian Reinforcement Learning with Regret Bounds](#)

- Brendan O'Donoghue
- abstract@[open-review](#): We consider the exploration-exploitation trade-off in reinforcement learning and show that an agent endowed with an exponential epistemic-risk-seeking utility function explores efficiently, as measured by regret. The state-action values induced by the exponential utility satisfy a Bellman recursion, so we can use dynamic programming to compute them. We call the resulting algorithm K-learning (for knowledge) and the risk-seeking utility ensures that the associated state-action values (K-values) are optimistic for the expected optimal Q-values under the posterior. The exponential utility function induces a Boltzmann exploration policy for which the 'temperature' parameter is equal to the risk-seeking parameter and is carefully controlled to yield a Bayes regret bound of $\tilde{O}(L^{3/2} \sqrt{SAT})$, where L is the time horizon, S is the number of states, A is the number of actions, and T is the total number of elapsed timesteps. We conclude with a numerical example demonstrating that K-learning is competitive with other state-of-the-art algorithms in practice.

[Logarithmic Regret from Sublinear Hints](#)

- Aditya Bhaskara, Ashok Cutkosky, Ravi Kumar, Manish Purohit
- abstract@[open-review](#): We consider the online linear optimization problem, where at every step the algorithm plays a point x_t in the unit ball, and suffers loss $\langle c_t, x_t \rangle$ for some cost vector c_t that is then revealed to the algorithm. Recent work showed that if an algorithm receives a hint h_t that has non-trivial correlation with c_t before it plays x_t , then it can achieve a regret guarantee of $O(\log T)$, improving on the bound

of $\Theta(\sqrt{T})$ in the standard setting. In this work, we study the question of whether an algorithm really requires a hint at *every* time step. Somewhat surprisingly, we show that an algorithm can obtain $O(\log T)$ regret with just $O(\sqrt{T})$ hints under a natural query model; in contrast, we also show that $O(\sqrt{T})$ hints cannot guarantee better than $\Omega(\sqrt{T})$ regret. We give two applications of our result, to the well-studied setting of $\{\text{em}\text{ optimistic}\}$ regret bounds, and to the problem of online learning with abstention.

[Independent mechanism analysis, a new concept?](#)

- Luigi Gresele, Julius von Kügelgen, Vincent Stimper, Bernhard Schölkopf, Michel Besserve
- abstract@[open-review](#): Independent component analysis provides a principled framework for unsupervised representation learning, with solid theory on the identifiability of the latent code that generated the data, given only observations of mixtures thereof. Unfortunately, when the mixing is nonlinear, the model is provably nonidentifiable, since statistical independence alone does not sufficiently constrain the problem. Identifiability can be recovered in settings where additional, typically observed variables are included in the generative process. We investigate an alternative path and consider instead including assumptions reflecting the principle of independent causal mechanisms exploited in the field of causality. Specifically, our approach is motivated by thinking of each source as independently influencing the mixing process. This gives rise to a framework which we term independent mechanism analysis. We provide theoretical and empirical evidence that our approach circumvents a number of nonidentifiability issues arising in nonlinear blind source separation.

[Momentum Centering and Asynchronous Update for Adaptive Gradient Methods](#)

- Juntang Zhuang, Yifan Ding, Tommy Tang, Nicha Dvornek, Sekhar C Tatikonda, James Duncan
- abstract@[open-review](#): We propose ACProp (Asynchronous-centering-Prop), an adaptive optimizer which combines centering of second momentum and asynchronous update (e.g. for t -th update, denominator uses information up to step $t-1$, while numerator uses gradient at t -th step). ACProp has both strong theoretical properties and empirical performance. With the example by Reddi et al. (2018), we show that asynchronous optimizers (e.g. AdaShift, ACProp) have weaker convergence condition than synchronous optimizers (e.g. Adam, RMSProp, AdaBelief); within asynchronous optimizers, we show that centering of second momentum further weakens the convergence condition. We demonstrate that ACProp has a convergence rate of $O(\frac{1}{\sqrt{T}})$ for the stochastic non-convex case, which matches the oracle rate and outperforms the $O(\frac{\log T}{\sqrt{T}})$ rate of RMSProp and Adam. We validate ACProp in extensive empirical studies: ACProp outperforms both SGD and other adaptive optimizers in image classification with CNN, and outperforms well-tuned adaptive optimizers in the training of various GAN models, reinforcement learning and transformers. To sum up, ACProp has good theoretical properties including weak convergence condition and optimal convergence rate, and strong empirical performance including good generalization like SGD and training stability like Adam. We provide the implementation at \url{https://github.com/juntang-zhuang/ACProp-Optimizer}.

[Robustness via Uncertainty-aware Cycle Consistency](#)

- Uddesha Upadhyay, Yanbei Chen, Zeynep Akata
- abstract@[open-review](#): Unpaired image-to-image translation refers to learning inter-image-domain mapping without corresponding image pairs. Existing methods learn deterministic mappings without explicitly modelling the robustness to outliers or predictive uncertainty, leading to performance degradation when encountering unseen perturbations at test time. To address this, we propose a novel probabilistic method based on Uncertainty-aware Generalized Adaptive Cycle Consistency (UGAC), which models the per-pixel residual by generalized Gaussian distribution, capable of modelling heavy-tailed distributions. We compare our model with a wide variety of state-of-the-art methods on various challenging tasks including unpaired image translation of natural images spanning autonomous driving, maps, facades, and also in the medical imaging domain consisting of MRI. Experimental results demonstrate that our method exhibits stronger robustness towards unseen perturbations in test data. Code is released here: <https://github.com/ExplainableML/UncertaintyAwareCycleConsistency>.

[CBP: backpropagation with constraint on weight precision using a pseudo-Lagrange multiplier method](#)

- Guhyun Kim, Doo Seok Jeong
- abstract@[open-review](#): Backward propagation of errors (backpropagation) is a method to minimize objective functions (e.g., loss functions) of deep neural networks by identifying optimal sets of weights and biases. Imposing constraints on weight precision is often required to alleviate prohibitive workloads on hardware. Despite the remarkable success of backpropagation, the algorithm itself is not capable of considering such constraints unless additional algorithms are applied simultaneously. To address this issue, we propose the constrained backpropagation (CBP) algorithm based on the pseudo-Lagrange multiplier method to obtain the optimal set of weights that satisfy a given set of constraints. The defining characteristic of the proposed CBP algorithm is the utilization of a Lagrangian function (loss function plus constraint function) as its objective function. We considered various types of constraints — binary, ternary, one-bit shift, and two-bit shift weight constraints. As a post-training method, CBP applied to AlexNet, ResNet-18, ResNet-50, and GoogLeNet on ImageNet, which were pre-trained using the conventional backpropagation. For most cases, the proposed algorithm outperforms the state-of-the-art methods on ImageNet, e.g., 66.6%, 74.4%, and 64.0% top-1 accuracy for ResNet-18, ResNet-50, and GoogLeNet with binary weights, respectively. This highlights CBP as a learning algorithm to address diverse constraints with the minimal performance loss by employing appropriate constraint functions. The code for CBP is publicly available at \url{https://github.com/dooseokjeong/CBP}.

[On the Sample Complexity of Privately Learning Axis-Aligned Rectangles](#)

- Menachem Sadigurschi, Uri Stemmer
- abstract@[open-review](#): We revisit the fundamental problem of learning Axis-Aligned-Rectangles over a finite grid $X^d \subseteq \mathbb{R}^d$ with differential privacy. Existing results show that the sample complexity of this problem is at most $\min\left\{ d \cdot \log|X|; d^{1.5} \cdot (\log|X|)^{1.5} \right\}$. That is, existing constructions either require sample complexity that grows linearly with $\log|X|$, or else it grows super linearly with the dimension d . We present a novel algorithm that reduces the sample complexity to only $\tilde{O}(d \cdot \log|X|^{1.5})$, attaining a dimensionality optimal dependency without requiring the sample complexity to grow with $\log|X|$. The technique used in order to attain this improvement involves the deletion of "exposed" data-points on the go, in a fashion designed to avoid the cost of the adaptive composition theorems. The core of this technique may be of individual interest, introducing a new method for constructing statistically-efficient private algorithms.

[Implicit Sparse Regularization: The Impact of Depth and Early Stopping](#)

- Jiangyuan Li, Thanh Nguyen, Chinmay Hegde, Ka Wai Wong
- abstract@[open-review](#): In this paper, we study the implicit bias of gradient descent for sparse regression. We extend results on regression with quadratic parametrization, which amounts to depth-2 diagonal linear networks, to more general depth-\$N\$ networks, under more realistic settings of noise and correlated designs. We show that early stopping is crucial for gradient descent to converge to a sparse model, a phenomenon that we call *implicit sparse regularization*. This result is in sharp contrast to known results for noiseless and uncorrelated-design cases. We characterize the impact of depth and early stopping and show that for a general depth parameter \$N\$, gradient descent with early stopping achieves minimax optimal sparse recovery with sufficiently small initialization w_0 and step size η . In particular, we show that increasing depth enlarges the scale of working initialization and the early-stopping window so that this implicit sparse regularization effect is more likely to take place.

Efficient Generalization with Distributionally Robust Learning

- Soumyadip Ghosh, Mark Squillante, Ebisa Wollega
- abstract@[open-review](#): Distributionally robust learning (DRL) is increasingly seen as a viable method to train machine learning models for improved model generalization. These min-max formulations, however, are more difficult to solve. We provide a new stochastic gradient descent algorithm to efficiently solve this DRL formulation. Our approach applies gradient descent to the outer minimization formulation and estimates the gradient of the inner maximization based on a sample average approximation. The latter uses a subset of the data sampled without replacement in each iteration, progressively increasing the subset size to ensure convergence. We rigorously establish convergence to a near-optimal solution under standard regularity assumptions and, for strongly convex losses, match the best known $\mathcal{O}(\epsilon^{-1})$ rate of convergence up to a known threshold. Empirical results demonstrate the significant benefits of our approach over previous work in improving learning for model generalization.

No-regret Online Learning over Riemannian Manifolds

- Xi Wang, Zhipeng Tu, Yiguang Hong, Yingyi Wu, Guodong Shi
- abstract@[open-review](#): We consider online optimization over Riemannian manifolds, where a learner attempts to minimize a sequence of time-varying loss functions defined on Riemannian manifolds. Though many Euclidean online convex optimization algorithms have been proven useful in a wide range of areas, less attention has been paid to their Riemannian counterparts. In this paper, we study Riemannian online gradient descent (R-OGD) on Hadamard manifolds for both geodesically convex and strongly geodesically convex loss functions, and Riemannian bandit algorithm (R-BAN) on Hadamard homogeneous manifolds for geodesically convex functions. We establish upper bounds on the regrets of the problem with respect to time horizon, manifold curvature, and manifold dimension. We also find a universal lower bound for the achievable regret by constructing an online convex optimization problem on Hadamard manifolds. All the obtained regret bounds match the corresponding results provided in Euclidean spaces. Finally, some numerical experiments validate our theoretical results.

Landmark-Guided Subgoal Generation in Hierarchical Reinforcement Learning

- Junsu Kim, Younggyo Seo, Jinwoo Shin
- abstract@[open-review](#): Goal-conditioned hierarchical reinforcement learning (HRL) has shown promising results for solving complex and long-horizon RL tasks. However, the action space of high-level policy in the goal-conditioned HRL is often large, so it results in poor exploration, leading to inefficiency in training. In this paper, we present Hierarchical reinforcement learning Guided by Landmarks (HIGL), a novel framework for training a high-level policy with a reduced action space guided by landmarks, i.e., promising states to explore. The key component of HIGL is twofold: (a) sampling landmarks that are informative for exploration and (b) encouraging the high level policy to generate a subgoal towards a selected landmark. For (a), we consider two criteria: coverage of the entire visited state space (i.e., dispersion of states) and novelty of states (i.e., prediction error of a state). For (b), we select a landmark as the very first landmark in the shortest path in a graph whose nodes are landmarks. Our experiments demonstrate that our framework outperforms prior-arts across a variety of control tasks, thanks to efficient exploration guided by landmarks.

Minimax Regret for Stochastic Shortest Path

- Alon Cohen, Yonathan Efroni, Yishay Mansour, Aviv Rosenberg
- abstract@[open-review](#): We study the Stochastic Shortest Path (SSP) problem in which an agent has to reach a goal state in minimum total expected cost. In the learning formulation of the problem, the agent has no prior knowledge about the costs and dynamics of the model. She repeatedly interacts with the model for K episodes, and has to minimize her regret. In this work we show that the minimax regret for this setting is $\widetilde{\mathcal{O}}(\sqrt{(B_{\star}^2 + B_{\star})|S||A|K})$ where B_{\star} is a bound on the expected cost of the optimal policy from any state, S is the state space, and A is the action space. This matches the $\Omega(\sqrt{B_{\star}^2 |S||A|K})$ lower bound of Rosenberg et al. [2020] for $B_{\star} \geq 1$, and improves their regret bound by a factor of $\sqrt{|S|}$. For $B_{\star} < 1$ we prove a matching lower bound of $\Omega(\sqrt{B_{\star}|S||A|K})$. Our algorithm is based on a novel reduction from SSP to finite-horizon MDPs. To that end, we provide an algorithm for the finite-horizon setting whose leading term in the regret depends polynomially on the expected cost of the optimal policy and only logarithmically on the horizon.

Parametrized Quantum Policies for Reinforcement Learning

- Sofiene Jerbi, Casper Gyurik, Simon Marshall, Hans Briegel, Vedran Dunjko
- abstract@[open-review](#): With the advent of real-world quantum computing, the idea that parametrized quantum computations can be used as hypothesis families in a quantum-classical machine learning system is gaining increasing traction. Such hybrid systems have already shown the potential to tackle real-world tasks in supervised and generative learning, and recent works have established their provable advantages in special artificial tasks. Yet, in the case of reinforcement learning, which is arguably most challenging and where learning boosts would be extremely valuable, no proposal has been successful in solving even standard benchmarking tasks, nor in showing a theoretical learning advantage over classical algorithms. In this work, we achieve both. We propose a hybrid quantum-classical reinforcement learning model using very few qubits, which we show can be effectively trained to solve several standard benchmarking environments. Moreover, we demonstrate, and formally prove, the ability of parametrized quantum circuits to solve certain learning tasks that are intractable to classical models, including current state-of-art deep neural networks, under the widely-believed classical hardness of the discrete logarithm problem.

On Pathologies in KL-Regularized Reinforcement Learning from Expert Demonstrations

- Tim G. J. Rudner, Cong Lu, Michael A Osborne, Yarin Gal, Yee Teh
- abstract@[open-review](#): KL-regularized reinforcement learning from expert demonstrations has proved successful in improving the sample efficiency of deep reinforcement learning algorithms, allowing them to be applied to challenging physical real-world tasks. However, we show that KL-regularized reinforcement learning with behavioral reference policies derived from expert demonstrations can suffer from pathological training dynamics that can lead to slow, unstable, and suboptimal online learning. We show empirically that the pathology occurs for commonly chosen behavioral policy classes and demonstrate its impact on sample efficiency and online policy performance. Finally, we show that the pathology can be remedied by non-parametric behavioral reference policies and that this allows KL-regularized reinforcement learning to significantly outperform state-of-the-art approaches on a variety of challenging locomotion and dexterous hand manipulation tasks.

Conditional Generation Using Polynomial Expansions

- Grigoris Chrysos, Markos Georgopoulos, Yannis Panagakis
- abstract@[open-review](#): Generative modeling has evolved to a notable field of machine learning. Deep polynomial neural networks (PNNs) have demonstrated impressive results in unsupervised image generation, where the task is to map an input vector (i.e., noise) to a synthesized image. However, the success of PNNs has not been replicated in conditional generation tasks, such as super-resolution. Existing PNNs focus on single-variable polynomial expansions which do not fare well to two-variable inputs, i.e., the noise variable and the conditional variable. In this work, we introduce a general framework, called CoPE, that enables a polynomial expansion of two input variables and captures their auto- and cross-correlations. We exhibit how CoPE can be trivially augmented to accept an arbitrary number of input variables. CoPE is evaluated in five tasks (class-conditional generation, inverse problems, edges-to-image translation, image-to-image translation, attribute-guided generation) involving eight datasets. The thorough evaluation suggests

that CoPE can be useful for tackling diverse conditional generation tasks. The source code of CoPE is available at <https://github.com/grigorisg9gr/polynomialnetsforconditionalgeneration>.

[Efficient constrained sampling via the mirror-Langevin algorithm](#)

- Kwangjun Ahn, Sinho Chewi
- abstract@[open-review](#): We propose a new discretization of the mirror-Langevin diffusion and give a crisp proof of its convergence. Our analysis uses relative convexity/smoothness and self-concordance, ideas which originated in convex optimization, together with a new result in optimal transport that generalizes the displacement convexity of the entropy. Unlike prior works, our result both (1) requires much weaker assumptions on the mirror map and the target distribution, and (2) has vanishing bias as the step size tends to zero. In particular, for the task of sampling from a log-concave distribution supported on a compact set, our theoretical results are significantly better than the existing guarantees.

[Adaptive Online Packing-guided Search for POMDPs](#)

- Chenyang Wu, Guoyu Yang, Zongzhang Zhang, Yang Yu, Dong Li, Wulong Liu, Jianye Hao
- abstract@[open-review](#): The partially observable Markov decision process (POMDP) provides a general framework for modeling an agent's decision process with state uncertainty, and online planning plays a pivotal role in solving it. A belief is a distribution of states representing state uncertainty. Methods for large-scale POMDP problems rely on the same idea of sampling both states and observations. That is, instead of exact belief updating, a collection of sampled states is used to approximate the belief; instead of considering all possible observations, only a set of sampled observations are considered. Inspired by this, we take one step further and propose an online planning algorithm, Adaptive Online Packing-guided Search (AdaOPS), to better approximate beliefs with adaptive particle filter technique and balance estimation bias and variance by fusing similar observation branches. Theoretically, our algorithm is guaranteed to find an $\$epsilon$ -optimal policy with a high probability given enough planning time under some mild assumptions. We evaluate our algorithm on several tricky POMDP domains, and it outperforms the state-of-the-art in all of them.

[Turing Completeness of Bounded-Precision Recurrent Neural Networks](#)

- Stephen Chung, Hava Siegelmann
- abstract@[open-review](#): Previous works have proved that recurrent neural networks (RNNs) are Turing-complete. However, in the proofs, the RNNs allow for neurons with unbounded precision, which is neither practical in implementation nor biologically plausible. To remove this assumption, we propose a dynamically growing memory module made of neurons of fixed precision. The memory module dynamically recruits new neurons when more memories are needed, and releases them when memories become irrelevant. We prove that a 54-neuron bounded-precision RNN with growing memory modules can simulate a Universal Turing Machine, with time complexity linear in the simulated machine's time and independent of the memory size. The result is extendable to various other stack-augmented RNNs. Furthermore, we analyze the Turing completeness of both unbounded-precision and bounded-precision RNNs, revisiting and extending the theoretical foundations of RNNs.

[End-to-end Multi-modal Video Temporal Grounding](#)

- Yi-Wen Chen, Yi-Hsuan Tsai, Ming-Hsuan Yang
- abstract@[open-review](#): We address the problem of text-guided video temporal grounding, which aims to identify the time interval of a certain event based on a natural language description. Different from most existing methods that only consider RGB images as visual features, we propose a multi-modal framework to extract complementary information from videos. Specifically, we adopt RGB images for appearance, optical flow for motion, and depth maps for image structure. While RGB images provide abundant visual cues of certain events, the performance may be affected by background clutters. Therefore, we use optical flow to focus on large motion and depth maps to infer the scene configuration when the action is related to objects recognizable with their shapes. To integrate the three modalities more effectively and enable inter-modal learning, we design a dynamic fusion scheme with transformers to model the interactions between modalities. Furthermore, we apply intra-modal self-supervised learning to enhance feature representations across videos for each modality, which also facilitates multi-modal learning. We conduct extensive experiments on the Charades-STA and ActivityNet Captions datasets, and show that the proposed method performs favorably against state-of-the-art approaches.

[How Powerful are Performance Predictors in Neural Architecture Search?](#)

- Colin White, Arber Zela, Robin Ru, Yang Liu, Frank Hutter
- abstract@[open-review](#): Early methods in the rapidly developing field of neural architecture search (NAS) required fully training thousands of neural networks. To reduce this extreme computational cost, dozens of techniques have since been proposed to predict the final performance of neural architectures. Despite the success of such performance prediction methods, it is not well-understood how different families of techniques compare to one another, due to the lack of an agreed-upon evaluation metric and optimization for different constraints on the initialization time and query time. In this work, we give the first large-scale study of performance predictors by analyzing 31 techniques ranging from learning curve extrapolation, to weight-sharing, to supervised learning, to zero-cost proxies. We test a number of correlation- and rank-based performance measures in a variety of settings, as well as the ability of each technique to speed up predictor-based NAS frameworks. Our results act as recommendations for the best predictors to use in different settings, and we show that certain families of predictors can be combined to achieve even better predictive power, opening up promising research directions. We release our code, featuring a library of 31 performance predictors.

[Stylized Dialogue Generation with Multi-Pass Dual Learning](#)

- Jinpeng Li, Yingce Xia, Rui Yan, Hongda Sun, Dongyan Zhao, Tie-Yan Liu
- abstract@[open-review](#): Stylized dialogue generation, which aims to generate a given-style response for an input context, plays a vital role in intelligent dialogue systems. Considering there is no parallel data between the contexts and the responses of target style S1, existing works mainly use back translation to generate stylized synthetic data for training, where the data about context, target style S1 and an intermediate style S0 is used. However, the interaction among these texts is not fully exploited, and the pseudo contexts are not adequately modeled. To overcome the above difficulties, we propose multi-pass dual learning (MPDL), which leverages the duality among the context, response of style S1 and response of style S_0. MPDL builds mappings among the above three domains, where the context should be reconstructed by the MPDL framework, and the reconstruction error is used as the training signal. To evaluate the quality of synthetic data, we also introduce discriminators that effectively measure how a pseudo sequence matches the specific domain, and the evaluation result is used as the weight for that data. Evaluation results indicate that our method obtains significant improvement over previous baselines.

[Entropy-based adaptive Hamiltonian Monte Carlo](#)

- Marcel Hirt, Michalis Titsias, Petros Dellaportas
- abstract@[open-review](#): Hamiltonian Monte Carlo (HMC) is a popular Markov Chain Monte Carlo (MCMC) algorithm to sample from an unnormalized probability distribution. A leapfrog integrator is commonly used to implement HMC in practice, but its performance can be sensitive to the choice of mass matrix used therein. We develop a gradient-based algorithm that allows for the adaptation of the mass matrix by encouraging the leapfrog integrator to have high acceptance rates while also exploring all dimensions jointly. In contrast to previous work that adapt the hyperparameters of HMC using some

form of expected squared jumping distance, the adaptation strategy suggested here aims to increase sampling efficiency by maximizing an approximation of the proposal entropy. We illustrate that using multiple gradients in the HMC proposal can be beneficial compared to a single gradient-step in Metropolis-adjusted Langevin proposals. Empirical evidence suggests that the adaptation method can outperform different versions of HMC schemes by adjusting the mass matrix to the geometry of the target distribution and by providing some control on the integration time.

[Continual World: A Robotic Benchmark For Continual Reinforcement Learning](#)

- Maciej Wołczyk, Michał Zajęc, Razvan Pascanu, Łukasz Kuciński, Piotr Miłoś
- abstract@[open-review](#): Continual learning (CL) --- the ability to continuously learn, building on previously acquired knowledge --- is a natural requirement for long-lived autonomous reinforcement learning (RL) agents. While building such agents, one needs to balance opposing desiderata, such as constraints on capacity and compute, the ability to not catastrophically forget, and to exhibit positive transfer on new tasks. Understanding the right trade-off is conceptually and computationally challenging, which we argue has led the community to overly focus on catastrophic forgetting. In response to these issues, we advocate for the need to prioritize forward transfer and propose Continual World, a benchmark consisting of realistic and meaningfully diverse robotic tasks built on top of Meta-World as a testbed. Following an in-depth empirical evaluation of existing CL methods, we pinpoint their limitations and highlight unique algorithmic challenges in the RL setting. Our benchmark aims to provide a meaningful and computationally inexpensive challenge for the community and thus help better understand the performance of existing and future solutions. Information about the benchmark, including the open-source code, is available at <https://sites.google.com/view/continualworld>.

[Towards Best-of-All-Worlds Online Learning with Feedback Graphs](#)

- Liad Erez, Tomer Koren
- abstract@[open-review](#): We study the online learning with feedback graphs framework introduced by Mannor and Shamir (2011), in which the feedback received by the online learner is specified by a graph $\$G\$$ over the available actions. We develop an algorithm that simultaneously achieves regret bounds of the form: $\$O(\sqrt{\theta(G) T})\$$ with adversarial losses; $\$O(\theta(G) \text{polylog}(T))\$$ with stochastic losses; and $\$O(\theta(G) \text{polylog}(T) + \sqrt{\theta(G) C})\$$ with stochastic losses subject to $\$C\$$ adversarial corruptions. Here, $\$\\theta(G)\$$ is the $\$\\text{clique-}\\text{covering-}\\text{number}\$$ of the graph $\$G\$$. Our algorithm is an instantiation of Follow-the-Regularized-Leader with a novel regularization that can be seen as a product of a Tsallis entropy component (inspired by Zimmert and Seldin (2019)) and a Shannon entropy component (analyzed in the corrupted stochastic case by Amir et al. (2020)), thus subtly interpolating between the two forms of entropies. One of our key technical contributions is in establishing the convexity of this regularizer and controlling its inverse Hessian, despite its complex product structure.

[ViTAE: Vision Transformer Advanced by Exploring Intrinsic Inductive Bias](#)

- Yufei Xu, Qiming ZHANG, Jing Zhang, Dacheng Tao
- abstract@[open-review](#): Transformers have shown great potential in various computer vision tasks owing to their strong capability in modeling long-range dependency using the self-attention mechanism. Nevertheless, vision transformers treat an image as 1D sequence of visual tokens, lacking an intrinsic inductive bias (IB) in modeling local visual structures and dealing with scale variance. Alternatively, they require large-scale training data and longer training schedules to learn the IB implicitly. In this paper, we propose a new Vision Transformer Advanced by Exploring intrinsic IB from convolutions, i.e., ViTAE. Technically, ViTAE has several spatial pyramid reduction modules to downsample and embed the input image into tokens with rich multi-scale context by using multiple convolutions with different dilation rates. In this way, it acquires an intrinsic scale invariance IB and is able to learn robust feature representation for objects at various scales. Moreover, in each transformer layer, ViTAE has a convolution block in parallel to the multi-head self-attention module, whose features are fused and fed into the feed-forward network. Consequently, it has the intrinsic locality IB and is able to learn local features and global dependencies collaboratively. Experiments on ImageNet as well as downstream tasks prove the superiority of ViTAE over the baseline transformer and concurrent works. Source code and pretrained models will be available at <https://github.com/Annbless/ViTAE>.

[Open Rule Induction](#)

- Wanyun Cui, Xingran Chen
- abstract@[open-review](#): Rules have a number of desirable properties. It is easy to understand, infer new knowledge, and communicate with other inference systems. One weakness of the previous rule induction systems is that they only find rules within a knowledge base (KB) and therefore cannot generalize to more open and complex real-world rules. Recently, the language model (LM)-based rule generation are proposed to enhance the expressive power of the rules. In this paper, we revisit the differences between KB-based rule induction and LM-based rule generation. We argue that, while KB-based methods inducted rules by discovering data commonalities, the current LM-based methods are learning rules from rules". This limits these methods to only produce "canned" rules whose patterns are constrained by the annotated rules, while discarding the rich expressive power of LMs for free text. Therefore, in this paper, we propose the open rule induction problem, which aims to induce open rules utilizing the knowledge in LMs. Besides, we propose the Orion (\underline{o} pen \underline{r} ule \underline{i} nduct \underline{on}) system to automatically mine open rules from LMs without supervision of annotated rules. We conducted extensive experiments to verify the quality and quantity of the inducted open rules. Surprisingly, when applying the open rules in downstream tasks (i.e. relation extraction), these automatically inducted rules even outperformed the manually annotated rules.

[Post-Contextual-Bandit Inference](#)

- Aurelien Bibaut, Maria Dimakopoulou, Nathan Kallus, Antoine Chambaz, Mark van der Laan
- abstract@[open-review](#): Contextual bandit algorithms are increasingly replacing non-adaptive A/B tests in e-commerce, healthcare, and policymaking because they can both improve outcomes for study participants and increase the chance of identifying good or even best policies. To support credible inference on novel interventions at the end of the study, nonetheless, we still want to construct valid confidence intervals on average treatment effects, subgroup effects, or value of new policies. The adaptive nature of the data collected by contextual bandit algorithms, however, makes this difficult: standard estimators are no longer asymptotically normally distributed and classic confidence intervals fail to provide correct coverage. While this has been addressed in non-contextual settings by using stabilized estimators, variance stabilized estimators in the contextual setting pose unique challenges that we tackle for the first time in this paper. We propose the Contextual Adaptive Doubly Robust (CADR) estimator, a novel estimator for policy value that is asymptotically normal under contextual adaptive data collection. The main technical challenge in constructing CADR is designing adaptive and consistent conditional standard deviation estimators for stabilization. Extensive numerical experiments using 57 OpenML datasets demonstrate that confidence intervals based on CADR uniquely provide correct coverage.

[Revisiting Discriminator in GAN Compression: A Generator-discriminator Cooperative Compression Scheme](#)

- Shaojie Li, Jie Wu, Xuefeng Xiao, Fei Chao, Xudong Mao, Rongrong Ji
- abstract@[open-review](#): Recently, a series of algorithms have been explored for GAN compression, which aims to reduce tremendous computational overhead and memory usages when deploying GANs on resource-constrained edge devices. However, most of the existing GAN compression work only focuses on how to compress the generator, while fails to take the discriminator into account. In this work, we revisit the role of discriminator in GAN compression and design a novel generator-discriminator cooperative compression scheme for GAN compression, termed GCC. Within GCC, a selective activation discriminator automatically selects and activates convolutional channels according to a local capacity constraint and a global coordination constraint, which help maintain the Nash equilibrium with the lightweight generator during the adversarial training and avoid mode collapse. The original generator and discriminator are also optimized from scratch, to play as a teacher model to progressively refine the pruned generator and the selective

activation discriminator. A novel online collaborative distillation scheme is designed to take full advantage of the intermediate feature of the teacher generator and discriminator to further boost the performance of the lightweight generator. Extensive experiments on various GAN-based generation tasks demonstrate the effectiveness and generalization of GCC. Among them, GCC contributes to reducing 80% computational costs while maintains comparable performance in image translation tasks.

[Asymptotically Exact Error Characterization of Offline Policy Evaluation with Misspecified Linear Models](#)

- Kohei Miyaguchi
- abstract@[open-review](#): We consider the problem of offline policy evaluation~(OPE) with Markov decision processes~(MDPs), where the goal is to estimate the utility of given decision-making policies based on static datasets. Recently, theoretical understanding of OPE has been rapidly advanced under (approximate) realizability assumptions, i.e., where the environments of interest are well approximated with the given hypothetical models. On the other hand, the OPE under unrealizability has not been well understood as much as in the realizable setting despite its importance in real-world applications. To address this issue, we study the behavior of a simple existing OPE method called the linear direct method~(DM) under the unrealizability. Consequently, we obtain an asymptotically exact characterization of the OPE error in a doubly robust form. Leveraging this result, we also establish the nonparametric consistency of the tile-coding estimators under quite mild assumptions.

[Topographic VAEs learn Equivariant Capsules](#)

- T. Anderson Keller, Max Welling
- abstract@[open-review](#): In this work we seek to bridge the concepts of topographic organization and equivariance in neural networks. To accomplish this, we introduce the Topographic VAE: a novel method for efficiently training deep generative models with topographically organized latent variables. We show that such a model indeed learns to organize its activations according to salient characteristics such as digit class, width, and style on MNIST. Furthermore, through topographic organization over time (i.e. temporal coherence), we demonstrate how predefined latent space transformation operators can be encouraged for observed transformed input sequences -- a primitive form of unsupervised learned equivariance. We demonstrate that this model successfully learns sets of approximately equivariant features (i.e. "capsules") directly from sequences and achieves higher likelihood on correspondingly transforming test sequences. Equivariance is verified quantitatively by measuring the approximate commutativity of the inference network and the sequence transformations. Finally, we demonstrate approximate equivariance to complex transformations, expanding upon the capabilities of existing group equivariant neural networks.

[MobILE: Model-Based Imitation Learning From Observation Alone](#)

- Rahul Kidambi, Jonathan Chang, Wen Sun
- abstract@[open-review](#): This paper studies Imitation Learning from Observations alone (ILFO) where the learner is presented with expert demonstrations that consist only of states visited by an expert (without access to actions taken by the expert). We present a provably efficient model-based framework MobILE to solve the ILFO problem. MobILE involves carefully trading off exploration against imitation - this is achieved by integrating the idea of optimism in the face of uncertainty into the distribution matching imitation learning (IL) framework. We provide a unified analysis for MobILE, and demonstrate that MobILE enjoys strong performance guarantees for classes of MDP dynamics that satisfy certain well studied notions of complexity. We also show that the ILFO problem is strictly harder than the standard IL problem by reducing ILFO to a multi-armed bandit problem indicating that exploration is necessary for solving ILFO efficiently. We complement these theoretical results with experimental simulations on benchmark OpenAI Gym tasks that indicate the efficacy of MobILE. Code for implementing the MobILE framework is available at <https://github.com/rahulkidambi/MobILE-NeurIPS2021>.

[Few-Round Learning for Federated Learning](#)

- Younghyun Park, Dong-Jun Han, Do-Yeon Kim, Jun Seo, Jaekyun Moon
- abstract@[open-review](#): In federated learning (FL), a number of distributed clients targeting the same task collaborate to train a single global model without sharing their data. The learning process typically starts from a randomly initialized or some pretrained model. In this paper, we aim at designing an initial model based on which an arbitrary group of clients can obtain a global model for its own purpose, within only a few rounds of FL. The key challenge here is that the downstream tasks for which the pretrained model will be used are generally unknown when the initial model is prepared. Our idea is to take a meta-learning approach to construct the initial model so that any group with a possibly unseen task can obtain a high-accuracy global model within only R rounds of FL. Our meta-learning itself could be done via federated learning among willing participants and is based on an episodic arrangement to mimic the R rounds of FL followed by inference in each episode. Extensive experimental results show that our method generalizes well for arbitrary groups of clients and provides large performance improvements given the same overall communication/computation resources, compared to other baselines relying on known pretraining methods.

[On Path Integration of Grid Cells: Group Representation and Isotropic Scaling](#)

- Ruiqi Gao, Jianwen Xie, Xue-Xin Wei, Song-Chun Zhu, Ying Nian Wu
- abstract@[open-review](#): Understanding how grid cells perform path integration calculations remains a fundamental problem. In this paper, we conduct theoretical analysis of a general representation model of path integration by grid cells, where the 2D self-position is encoded as a higher dimensional vector, and the 2D self-motion is represented by a general transformation of the vector. We identify two conditions on the transformation. One is a group representation condition that is necessary for path integration. The other is an isotropic scaling condition that ensures locally conformal embedding, so that the error in the vector representation translates conformally to the error in the 2D self-position. Then we investigate the simplest transformation, i.e., the linear transformation, uncover its explicit algebraic and geometric structure as matrix Lie group of rotation, and explore the connection between the isotropic scaling condition and a special class of hexagon grid patterns. Finally, with our optimization-based approach, we manage to learn hexagon grid patterns that share similar properties of the grid cells in the rodent brain. The learned model is capable of accurate long distance path integration. Code is available at <https://github.com/ruiqigao/grid-cell-path>.

[Online Convex Optimization with Continuous Switching Constraint](#)

- Guanghui Wang, Yuanyu Wan, Tianbao Yang, Lijun Zhang
- abstract@[open-review](#): In many sequential decision making applications, the change of decision would bring an additional cost, such as the wear-and-tear cost associated with changing server status. To control the switching cost, we introduce the problem of online convex optimization with continuous switching constraint, where the goal is to achieve a small regret given a budget on the \emph{overall} switching cost. We first investigate the hardness of the problem, and provide a lower bound of order $\Omega(\sqrt{T})$ when the switching cost budget $S = \Omega(\sqrt{T})$, and $\Omega(\min\{\frac{T}{S}, T\})$ when $S = O(\sqrt{T})$, where T is the time horizon. The essential idea is to carefully design an adaptive adversary, who can adjust the loss function according to the cumulative switching cost of the player incurred so far based on the orthogonal technique. We then develop a simple gradient-based algorithm which enjoys the minimax optimal regret bound. Finally, we show that, for strongly convex functions, the regret bound can be improved to $O(\log T)$ for $S = \Omega(\log T)$, and $O(\min\{T \exp(S) + S, T\})$ for $S = O(\log T)$.

[Why Do Better Loss Functions Lead to Less Transferable Features?](#)

- Simon Kornblith, Ting Chen, Honglak Lee, Mohammad Norouzi
- abstract@[open-review](#): Previous work has proposed many new loss functions and regularizers that improve test accuracy on image classification tasks. However, it is not clear whether these loss functions learn better representations for downstream tasks. This paper studies how the choice of training objective affects the transferability of the hidden representations of convolutional neural networks trained on ImageNet. We show that many objectives lead to statistically significant improvements in ImageNet accuracy over vanilla softmax cross-entropy, but the resulting fixed feature extractors transfer substantially worse to downstream tasks, and the choice of loss has little effect when networks are fully fine-tuned on the new tasks. Using centered kernel alignment to measure similarity between hidden representations of networks, we find that differences among loss functions are apparent only in the last few layers of the network. We delve deeper into representations of the penultimate layer, finding that different objectives and hyperparameter combinations lead to dramatically different levels of class separation. Representations with higher class separation obtain higher accuracy on the original task, but their features are less useful for downstream tasks. Our results suggest there exists a trade-off between learning invariant features for the original task and features relevant for transfer tasks.

[**Breaking the centralized barrier for cross-device federated learning**](#)

- Sai Praneeth Karimireddy, Martin Jaggi, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian U. Stich, Ananda Theertha Suresh
- abstract@[open-review](#): Federated learning (FL) is a challenging setting for optimization due to the heterogeneity of the data across different clients which gives rise to the client drift phenomenon. In fact, obtaining an algorithm for FL which is uniformly better than simple centralized training has been a major open problem thus far. In this work, we propose a general algorithmic framework, Mime, which i) mitigates client drift and ii) adapts arbitrary centralized optimization algorithms such as momentum and Adam to the cross-device federated learning setting. Mime uses a combination of control-variates and server-level statistics (e.g. momentum) at every client-update step to ensure that each local update mimics that of the centralized method run on iid data. We prove a reduction result showing that Mime can translate the convergence of a generic algorithm in the centralized setting into convergence in the federated setting. Further, we show that when combined with momentum based variance reduction, Mime is provably faster than any centralized method--the first such result. We also perform a thorough experimental exploration of Mime's performance on real world datasets.

[**Adversarially robust learning for security-constrained optimal power flow**](#)

- Priya Donti, Aayushya Agarwal, Neeraj Vijay Bedmutha, Larry Pileggi, J. Zico Kolter
- abstract@[open-review](#): In recent years, the ML community has seen surges of interest in both adversarially robust learning and implicit layers, but connections between these two areas have seldom been explored. In this work, we combine innovations from these areas to tackle the problem of N-k security-constrained optimal power flow (SCOPF). N-k SCOPF is a core problem for the operation of electrical grids, and aims to schedule power generation in a manner that is robust to potentially \$k\$ simultaneous equipment outages. Inspired by methods in adversarially robust training, we frame N-k SCOPF as a minimax optimization problem -- viewing power generation settings as adjustable parameters and equipment outages as (adversarial) attacks -- and solve this problem via gradient-based techniques. The loss function of this minimax problem involves resolving implicit equations representing grid physics and operational decisions, which we differentiate through via the implicit function theorem. We demonstrate the efficacy of our framework in solving N-3 SCOPF, which has traditionally been considered as prohibitively expensive to solve given that the problem size depends combinatorially on the number of potential outages.

[**Learning a Single Neuron with Bias Using Gradient Descent**](#)

- Gal Vardi, Gilad Yehudai, Ohad Shamir
- abstract@[open-review](#): We theoretically study the fundamental problem of learning a single neuron with a bias term ($\mathbf{x} \mapsto \sigma(\langle \mathbf{w}, \mathbf{x} \rangle + b)$) in the realizable setting with the ReLU activation, using gradient descent. Perhaps surprisingly, we show that this is a significantly different and more challenging problem than the bias-less case (which was the focus of previous works on single neurons), both in terms of the optimization geometry as well as the ability of gradient methods to succeed in some scenarios. We provide a detailed study of this problem, characterizing the critical points of the objective, demonstrating failure cases, and providing positive convergence guarantees under different sets of assumptions. To prove our results, we develop some tools which may be of independent interest, and improve previous results on learning single neurons.

[**Making a \(Counterfactual\) Difference One Rationale at a Time**](#)

- Mitchell Plyler, Michael Green, Min Chi
- abstract@[open-review](#): Rationales, snippets of extracted text that explain an inference, have emerged as a popular framework for interpretable natural language processing (NLP). Rationale models typically consist of two cooperating modules: a selector and a classifier with the goal of maximizing the mutual information (MMI) between the "selected" text and the document label. Despite their promises, MMI-based methods often pick up on spurious text patterns and result in models with nonsensical behaviors. In this work, we investigate whether counterfactual data augmentation (CDA), without human assistance, can improve the performance of the selector by lowering the mutual information between spurious signals and the document label. Our counterfactuals are produced in an unsupervised fashion using class-dependent generative models. From an information theoretic lens, we derive properties of the unaugmented dataset for which our CDA approach would succeed. The effectiveness of CDA is empirically evaluated by comparing against several baselines including an improved MMI-based rationale schema on two multi-aspect datasets. Our results show that CDA produces rationales that better capture the signal of interest.

[**3D Siamese Voxel-to-BEV Tracker for Sparse Point Clouds**](#)

- Le Hui, Lingpeng Wang, Mingmei Cheng, Jin Xie, Jian Yang
- abstract@[open-review](#): 3D object tracking in point clouds is still a challenging problem due to the sparsity of LiDAR points in dynamic environments. In this work, we propose a Siamese voxel-to-BEV tracker, which can significantly improve the tracking performance in sparse 3D point clouds. Specifically, it consists of a Siamese shape-aware feature learning network and a voxel-to-BEV target localization network. The Siamese shape-aware feature learning network can capture 3D shape information of the object to learn the discriminative features of the object so that the potential target from the background in sparse point clouds can be identified. To this end, we first perform template feature embedding to embed the template's feature into the potential target and then generate a dense 3D shape to characterize the shape information of the potential target. For localizing the tracked target, the voxel-to-BEV target localization network regresses the target's 2D center and the z-axis center from the dense bird's eye view (BEV) feature map in an anchor-free manner. Concretely, we compress the voxelized point cloud along z-axis through max pooling to obtain a dense BEV feature map, where the regression of the 2D center and the z-axis center can be performed more effectively. Extensive evaluation on the KITTI tracking dataset shows that our method significantly outperforms the current state-of-the-art methods by a large margin. Code is available at <https://github.com/fptthink/V2B>.

[**Stateful Strategic Regression**](#)

- Keegan Harris, Hoda Heidari, Steven Z. Wu
- abstract@[open-review](#): Automated decision-making tools increasingly assess individuals to determine if they qualify for high-stakes opportunities. A recent line of research investigates how strategic agents may respond to such scoring tools to receive favorable assessments. While prior work has focused on the short-term strategic interactions between a decision-making institution (modeled as a principal) and individual decision-subjects (modeled as agents), we investigate interactions spanning multiple time-steps. In particular, we consider settings in which the agent's effort investment today can

accumulate over time in the form of an internal state - impacting both his future rewards and that of the principal. We characterize the Stackelberg equilibrium of the resulting game and provide novel algorithms for computing it. Our analysis reveals several intriguing insights about the role of multiple interactions in shaping the game's outcome: First, we establish that in our stateful setting, the class of all linear assessment policies remains as powerful as the larger class of all monotonic assessment policies. While recovering the principal's optimal policy requires solving a non-convex optimization problem, we provide polynomial-time algorithms for recovering both the principal and agent's optimal policies under common assumptions about the process by which effort investments convert to observable features. Most importantly, we show that with multiple rounds of interaction at her disposal, the principal is more effective at incentivizing the agent to accumulate effort in her desired direction. Our work addresses several critical gaps in the growing literature on the societal impacts of automated decision-making - by focusing on longer time horizons and accounting for the compounding nature of decisions individuals receive over time.

[Self-Attention Between Datapoints: Going Beyond Individual Input-Output Pairs in Deep Learning](#)

- Jannik Kossen, Neil Band, Clare Lyle, Aidan N. Gomez, Thomas Rainforth, Yarin Gal
- abstract@[open-review](#): We challenge a common assumption underlying most supervised deep learning: that a model makes a prediction depending only on its parameters and the features of a single input. To this end, we introduce a general-purpose deep learning architecture that takes as input the entire dataset instead of processing one datapoint at a time. Our approach uses self-attention to reason about relationships between datapoints explicitly, which can be seen as realizing non-parametric models using parametric attention mechanisms. However, unlike conventional non-parametric models, we let the model learn end-to-end from the data how to make use of other datapoints for prediction. Empirically, our models solve cross-datapoint lookup and complex reasoning tasks unsolvable by traditional deep learning models. We show highly competitive results on tabular data, early results on CIFAR-10, and give insight into how the model makes use of the interactions between points.

[Your head is there to move you around: Goal-driven models of the primate dorsal pathway](#)

- Patrick Mineault, Shahab Bakhtiari, Blake Richards, Christopher Pack
- abstract@[open-review](#): Neurons in the dorsal visual pathway of the mammalian brain are selective for motion stimuli, with the complexity of stimulus representations increasing along the hierarchy. This progression is similar to that of the ventral visual pathway, which is well characterized by artificial neural networks (ANNs) optimized for object recognition. In contrast, there are no image-computable models of the dorsal stream with comparable explanatory power. We hypothesized that the properties of dorsal stream neurons could be explained by a simple learning objective: the need for an organism to orient itself during self-motion. To test this hypothesis, we trained a 3D ResNet to predict an agent's self-motion parameters from visual stimuli in a simulated environment. We found that the responses in this network accounted well for the selectivity of neurons in a large database of single-neuron recordings from the dorsal visual stream of non-human primates. In contrast, ANNs trained on an action recognition dataset through supervised or self-supervised learning could not explain responses in the dorsal stream, despite also being trained on naturalistic videos with moving objects. These results demonstrate that an ecologically relevant cost function can account for dorsal stream properties in the primate brain.

[Achieving Rotational Invariance with Bessel-Convolutional Neural Networks](#)

- Valentin Delchevalerie, Adrien Bibal, Benoît Frénay, Alexandre Mayer
- abstract@[open-review](#): For many applications in image analysis, learning models that are invariant to translations and rotations is paramount. This is the case, for example, in medical imaging where the objects of interest can appear at arbitrary positions, with arbitrary orientations. As of today, Convolutional Neural Networks (CNN) are one of the most powerful tools for image analysis. They achieve, thanks to convolutions, an invariance with respect to translations. In this work, we present a new type of convolutional layer that takes advantage of Bessel functions, well known in physics, to build Bessel-CNNs (B-CNNs) that are invariant to all the continuous set of possible rotation angles by design.

[Unsupervised Domain Adaptation with Dynamics-Aware Rewards in Reinforcement Learning](#)

- Jinxin Liu, Hao Shen, Donglin Wang, Yuchen Kang, Qiangxing Tian
- abstract@[open-review](#): Unsupervised reinforcement learning aims to acquire skills without prior goal representations, where an agent automatically explores an open-ended environment to represent goals and learn the goal-conditioned policy. However, this procedure is often time-consuming, limiting the rollout in some potentially expensive target environments. The intuitive approach of training in another interaction-rich environment disrupts the reproducibility of trained skills in the target environment due to the dynamics shifts and thus inhibits direct transferring. Assuming free access to a source environment, we propose an unsupervised domain adaptation method to identify and acquire skills across dynamics. Particularly, we introduce a KL regularized objective to encourage emergence of skills, rewarding the agent for both discovering skills and aligning its behaviors respecting dynamics shifts. This suggests that both dynamics (source and target) shape the reward to facilitate the learning of adaptive skills. We also conduct empirical experiments to demonstrate that our method can effectively learn skills that can be smoothly deployed in target.

[GraphFormers: GNN-nested Transformers for Representation Learning on Textual Graph](#)

- Junhan Yang, Zheng Liu, Shitao Xiao, Chaozhuo Li, Defu Lian, Sanjay Agrawal, Amit Singh, Guangzhong Sun, Xing Xie
- abstract@[open-review](#): The representation learning on textual graph is to generate low-dimensional embeddings for the nodes based on the individual textual features and the neighbourhood information. Recent breakthroughs on pretrained language models and graph neural networks push forward the development of corresponding techniques. The existing works mainly rely on the cascaded model architecture: the textual features of nodes are independently encoded by language models at first; the textual embeddings are aggregated by graph neural networks afterwards. However, the above architecture is limited due to the independent modeling of textual features. In this work, we propose GraphFormers, where layerwise GNN components are nested alongside the transformer blocks of language models. With the proposed architecture, the text encoding and the graph aggregation are fused into an iterative workflow, making each node's semantic accurately comprehended from the global perspective. In addition, a progressive learning strategy is introduced, where the model is successively trained on manipulated data and original data to reinforce its capability of integrating information on graph. Extensive evaluations are conducted on three large-scale benchmark datasets, where GraphFormers outperform the SOTA baselines with comparable running efficiency. The source code is released at <https://github.com/microsoft/GraphFormers>.

[A Universal Law of Robustness via Isoperimetry](#)

- Sébastien Bubeck, Mark Sellke
- abstract@[open-review](#): Classically, data interpolation with a parametrized model class is possible as long as the number of parameters is larger than the number of equations to be satisfied. A puzzling phenomenon in the current practice of deep learning is that models are trained with many more parameters than what this classical theory would suggest. We propose a theoretical explanation for this phenomenon. We prove that for a broad class of data distributions and model classes, overparametrization is necessary if one wants to interpolate the data smoothly. Namely we show that smooth interpolation requires d times more parameters than mere interpolation, where d is the ambient data dimension. We prove this universal law of robustness for any smoothly parametrized function class with polynomial size weights, and any covariate distribution verifying isoperimetry. In the case of two-layers neural networks and Gaussian covariates, this law was conjectured in prior work by Bubeck, Li and Nagaraj. We also give an interpretation of our result as an improved generalization bound for model classes consisting of smooth functions.

[On Contrastive Representations of Stochastic Processes](#)

- Emile Mathieu, Adam Foster, Yee Teh
- abstract@[open-review](#): Learning representations of stochastic processes is an emerging problem in machine learning with applications from meta-learning to physical object models to time series. Typical methods rely on exact reconstruction of observations, but this approach breaks down as observations become high-dimensional or noise distributions become complex. To address this, we propose a unifying framework for learning contrastive representations of stochastic processes (CReSP) that does away with exact reconstruction. We dissect potential use cases for stochastic process representations, and propose methods that accommodate each. Empirically, we show that our methods are effective for learning representations of periodic functions, 3D objects and dynamical processes. Our methods tolerate noisy high-dimensional observations better than traditional approaches, and the learned representations transfer to a range of downstream tasks.

[A Domain-Shrinking based Bayesian Optimization Algorithm with Order-Optimal Regret Performance](#)

- Sudeep Salgia, Sattar Vakili, Qing Zhao
- abstract@[open-review](#): We consider sequential optimization of an unknown function in a reproducing kernel Hilbert space. We propose a Gaussian process-based algorithm and establish its order-optimal regret performance (up to a poly-logarithmic factor). This is the first GP-based algorithm with an order-optimal regret guarantee. The proposed algorithm is rooted in the methodology of domain shrinking realized through a sequence of tree-based region pruning and refining to concentrate queries in increasingly smaller high-performing regions of the function domain. The search for high-performing regions is localized and guided by an iterative estimation of the optimal function value to ensure both learning efficiency and computational efficiency. Compared with the prevailing GP-UCB family of algorithms, the proposed algorithm reduces computational complexity by a factor of $\mathcal{O}(T^{2d-1})$ (where T is the time horizon and d the dimension of the function domain).

[Scalars are universal: Equivariant machine learning, structured like classical physics](#)

- Soledad Villar, David W Hogg, Kate Storey-Fisher, Weichi Yao, Ben Blum-Smith
- abstract@[open-review](#): There has been enormous progress in the last few years in designing neural networks that respect the fundamental symmetries and coordinate freedoms of physical law. Some of these frameworks make use of irreducible representations, some make use of high-order tensor objects, and some apply symmetry-enforcing constraints. Different physical laws obey different combinations of fundamental symmetries, but a large fraction (possibly all) of classical physics is equivariant to translation, rotation, reflection (parity), boost (relativity), and permutations. Here we show that it is simple to parameterize universally approximating polynomial functions that are equivariant under these symmetries, or under the Euclidean, Lorentz, and Poincaré groups, at any dimensionality d . The key observation is that nonlinear $O(d)$ -equivariant (and related-group-equivariant) functions can be universally expressed in terms of a lightweight collection of scalars---scalar products and scalar contractions of the scalar, vector, and tensor inputs. We complement our theory with numerical examples that show that the scalar-based method is simple, efficient, and scalable.

[Unsupervised Object-Level Representation Learning from Scene Images](#)

- Jiahao Xie, Xiaohang Zhan, Ziwei Liu, Yew Soon Ong, Chen Change Loy
- abstract@[open-review](#): Contrastive self-supervised learning has largely narrowed the gap to supervised pre-training on ImageNet. However, its success highly relies on the object-centric priors of ImageNet, i.e., different augmented views of the same image correspond to the same object. Such a heavily curated constraint becomes immediately infeasible when pre-trained on more complex scene images with many objects. To overcome this limitation, we introduce Object-level Representation Learning (ORL), a new self-supervised learning framework towards scene images. Our key insight is to leverage image-level self-supervised pre-training as the prior to discover object-level semantic correspondence, thus realizing object-level representation learning from scene images. Extensive experiments on COCO show that ORL significantly improves the performance of self-supervised learning on scene images, even surpassing supervised ImageNet pre-training on several downstream tasks. Furthermore, ORL improves the downstream performance when more unlabeled scene images are available, demonstrating its great potential of harnessing unlabeled data in the wild. We hope our approach can motivate future research on more general-purpose unsupervised representation learning from scene data.

[Do Transformers Really Perform Badly for Graph Representation?](#)

- Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, Tie-Yan Liu
- abstract@[open-review](#): The Transformer architecture has become a dominant choice in many domains, such as natural language processing and computer vision. Yet, it has not achieved competitive performance on popular leaderboards of graph-level prediction compared to mainstream GNN variants. Therefore, it remains a mystery how Transformers could perform well for graph representation learning. In this paper, we solve this mystery by presenting Graphomer, which is built upon the standard Transformer architecture, and could attain excellent results on a broad range of graph representation learning tasks, especially on the recent OGB Large-Scale Challenge. Our key insight to utilizing Transformer in the graph is the necessity of effectively encoding the structural information of a graph into the model. To this end, we propose several simple yet effective structural encoding methods to help Graphomer better model graph-structured data. Besides, we mathematically characterize the expressive power of Graphomer and exhibit that with our ways of encoding the structural information of graphs, many popular GNN variants could be covered as the special cases of Graphomer. The code and models of Graphomer will be made publicly available at \url{https://github.com/Microsoft/Graphomer}.

[Powerpropagation: A sparsity inducing weight reparameterisation](#)

- Jonathan Schwarz, Siddhant Jayakumar, Razvan Pascanu, Peter E Latham, Yee Teh
- abstract@[open-review](#): The training of sparse neural networks is becoming an increasingly important tool for reducing the computational footprint of models at training and evaluation, as well enabling the effective scaling up of models. Whereas much work over the years has been dedicated to specialised pruning techniques, little attention has been paid to the inherent effect of gradient based training on model sparsity. In this work, we introduce Powerpropagation, a new weight-parameterisation for neural networks that leads to inherently sparse models. Exploiting the behaviour of gradient descent, our method gives rise to weight updates exhibiting a “rich get richer” dynamic, leaving low-magnitude parameters largely unaffected by learning. Models trained in this manner exhibit similar performance, but have a distribution with markedly higher density at zero, allowing more parameters to be pruned safely. Powerpropagation is general, intuitive, cheap and straight-forward to implement and can readily be combined with various other techniques. To highlight its versatility, we explore it in two very different settings: Firstly, following a recent line of work, we investigate its effect on sparse training for resource-constrained settings. Here, we combine Powerpropagation with a traditional weight-pruning technique as well as recent state-of-the-art sparse-to-sparse algorithms, showing superior performance on the ImageNet benchmark. Secondly, we advocate the use of sparsity in overcoming catastrophic forgetting, where compressed representations allow accommodating a large number of tasks at fixed model capacity. In all cases our reparameterisation considerably increases the efficacy of the off-the-shelf methods.

[Stronger NAS with Weaker Predictors](#)

- Junru Wu, Xiyang Dai, Dongdong Chen, Yinpeng Chen, Mengchen Liu, Ye Yu, Zhangyang Wang, Zicheng Liu, Mei Chen, Lu Yuan
- abstract@[open-review](#): Neural Architecture Search (NAS) often trains and evaluates a large number of architectures. Recent predictor-based NAS approaches attempt to alleviate such heavy computation costs with two key steps: sampling some architecture-performance pairs and fitting a proxy

accuracy predictor. Given limited samples, these predictors, however, are far from accurate to locate top architectures due to the difficulty of fitting the huge search space. This paper reflects on a simple yet crucial question: if our final goal is to find the best architecture, do we really need to model the whole space well?. We propose a paradigm shift from fitting the whole architecture space using one strong predictor, to progressively fitting a search path towards the high-performance sub-space through a set of weaker predictors. As a key property of the weak predictors, their probabilities of sampling better architectures keep increasing. Hence we only sample a few well-performed architectures guided by the previously learned predictor and estimate a new better weak predictor. This embarrassingly easy framework, dubbed WeakNAS, produces coarse-to-fine iteration to gradually refine the ranking of sampling space. Extensive experiments demonstrate that WeakNAS costs fewer samples to find top-performance architectures on NAS-Bench-101 and NAS-Bench-201. Compared to state-of-the-art (SOTA) predictor-based NAS methods, WeakNAS outperforms all with notable margins, e.g., requiring at least 7.5x less samples to find global optimal on NAS-Bench-101. WeakNAS can also absorb their ideas to boost performance more. Further, WeakNAS strikes the new SOTA result of 81.3% in the ImageNet MobileNet Search Space. The code is available at: <https://github.com/VITA-Group/WeakNAS>.

[Convolutional Normalization: Improving Deep Convolutional Network Robustness and Training](#)

- Sheng Liu, Xiao Li, Yueyang Zhai, Chong You, Zhihui Zhu, Carlos Fernandez-Granda, Qing Qu
- abstract@[open-review](#): Normalization techniques have become a basic component in modern convolutional neural networks (ConvNets). In particular, many recent works demonstrate that promoting the orthogonality of the weights helps train deep models and improve robustness. For ConvNets, most existing methods are based on penalizing or normalizing weight matrices derived from concatenating or flattening the convolutional kernels. These methods often destroy or ignore the benign convolutional structure of the kernels; therefore, they are often expensive or impractical for deep ConvNets. In contrast, we introduce a simple and efficient ``Convolutional Normalization'' (ConvNorm) method that can fully exploit the convolutional structure in the Fourier domain and serve as a simple plug-and-play module to be conveniently incorporated into any ConvNets. Our method is inspired by recent work on preconditioning methods for convolutional sparse coding and can effectively promote each layer's channel-wise isometry. Furthermore, we show that our ConvNorm can reduce the layerwise spectral norm of the weight matrices and hence improve the Lipschitzness of the network, leading to easier training and improved robustness for deep ConvNets. Applied to classification under noise corruptions and generative adversarial network (GAN), we show that the ConvNorm improves the robustness of common ConvNets such as ResNet and the performance of GAN. We verify our findings via numerical experiments on CIFAR and ImageNet. Our implementation is available online at \url{<https://github.com/shengliu66/ConvNorm>}.

[Nearly-Tight and Oblivious Algorithms for Explainable Clustering](#)

- Buddhima Gamlath, Xinrui Jia, Adam Polak, Ola Svensson
- abstract@[open-review](#): We study the problem of explainable clustering in the setting first formalized by Dasgupta, Frost, Moshkovitz, and Rashtchian (ICML 2020). A $\$k\$$ -clustering is said to be explainable if it is given by a decision tree where each internal node splits data points with a threshold cut in a single dimension (feature), and each of the $\$k\$$ leaves corresponds to a cluster. We give an algorithm that outputs an explainable clustering that loses at most a factor of $\$O(\log^2 k)$ compared to an optimal (not necessarily explainable) clustering for the $\$k\$$ -medians objective, and a factor of $\$O(k \log^2 k)$ for the $\$k\$$ -means objective. This improves over the previous best upper bounds of $\$O(k)$ and $\$O(k^2)$, respectively, and nearly matches the previous $\$\\Omega(\log k)$ lower bound for $\$k\$$ -medians and our new $\$\\Omega(k)$ lower bound for $\$k\$$ -means. The algorithm is remarkably simple. In particular, given an initial not necessarily explainable clustering in $\$\\mathbb{R}^d$, it is oblivious to the data points and runs in time $\$O(dk \log^2 k)$, independent of the number of data points $\$n\$$. Our upper and lower bounds also generalize to objectives given by higher $\$\\ell_p$ -norms.

[Deep Networks Provably Classify Data on Curves](#)

- Tingran Wang, Sam Buchanan, Dar Gilboa, John Wright
- abstract@[open-review](#): Data with low-dimensional nonlinear structure are ubiquitous in engineering and scientific problems. We study a model problem with such structure---a binary classification task that uses a deep fully-connected neural network to classify data drawn from two disjoint smooth curves on the unit sphere. Aside from mild regularity conditions, we place no restrictions on the configuration of the curves. We prove that when (i) the network depth is large relative to certain geometric properties that set the difficulty of the problem and (ii) the network width and number of samples is polynomial in the depth, randomly-initialized gradient descent quickly learns to correctly classify all points on the two curves with high probability. To our knowledge, this is the first generalization guarantee for deep networks with nonlinear data that depends only on intrinsic data properties. Our analysis proceeds by a reduction to dynamics in the neural tangent kernel (NTK) regime, where the network depth plays the role of a fitting resource in solving the classification problem. In particular, via fine-grained control of the decay properties of the NTK, we demonstrate that when the network is sufficiently deep, the NTK can be locally approximated by a translationally invariant operator on the manifolds and stably inverted over smooth functions, which guarantees convergence and generalization.

[COMBO: Conservative Offline Model-Based Policy Optimization](#)

- Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, Chelsea Finn
- abstract@[open-review](#): Model-based reinforcement learning (RL) algorithms, which learn a dynamics model from logged experience and perform conservative planning under the learned model, have emerged as a promising paradigm for offline reinforcement learning (offline RL). However, practical variants of such model-based algorithms rely on explicit uncertainty quantification for incorporating conservatism. Uncertainty estimation with complex models, such as deep neural networks, can be difficult and unreliable. We empirically find that uncertainty estimation is not accurate and leads to poor performance in certain scenarios in offline model-based RL. We overcome this limitation by developing a new model-based offline RL algorithm, COMBO, that trains a value function using both the offline dataset and data generated using rollouts under the model while also additionally regularizing the value function on out-of-support state-action tuples generated via model rollouts. This results in a conservative estimate of the value function for out-of-support state-action tuples, without requiring explicit uncertainty estimation. Theoretically, we show that COMBO satisfies a policy improvement guarantee in the offline setting. Through extensive experiments, we find that COMBO attains greater performance compared to prior offline RL on problems that demand generalization to related but previously unseen tasks, and also consistently matches or outperforms prior offline RL methods on widely studied offline RL benchmarks, including image-based tasks.

[Time-series Generation by Contrastive Imitation](#)

- Daniel Jarrett, Ioana Bica, Mihaela van der Schaar
- abstract@[open-review](#): Consider learning a generative model for time-series data. The sequential setting poses a unique challenge: Not only should the generator capture the conditional dynamics of (stepwise) transitions, but its open-loop rollouts should also preserve the joint distribution of (multi-step) trajectories. On one hand, autoregressive models trained by MLE allow learning and computing explicit transition distributions, but suffer from compounding error during rollouts. On the other hand, adversarial models based on GAN training alleviate such exposure bias, but transitions are implicit and hard to assess. In this work, we study a generative framework that seeks to combine the strengths of both: Motivated by a moment-matching objective to mitigate compounding error, we optimize a local (but forward-looking) transition policy, where the reinforcement signal is provided by a global (but stepwise-decomposable) energy model trained by contrastive estimation. At training, the two components are learned cooperatively, avoiding the instabilities typical of adversarial objectives. At inference, the learned policy serves as the generator for iterative sampling, and the learned energy serves as a trajectory-level measure for evaluating sample quality. By expressly training a policy to imitate sequential behavior of time-series features in a dataset, this approach embodies "generation by imitation". Theoretically, we illustrate the correctness of this formulation and the consistency of the algorithm. Empirically, we evaluate its ability to generate predictively useful samples from real-world datasets, verifying that it performs at the standard of existing benchmarks.

Differentially Private Sampling from Distributions

- Sofya Raskhodnikova, Satchit Sivakumar, Adam Smith, Marika Swanberg
- abstract@[open-review](#): We initiate an investigation of private sampling from distributions. Given a dataset with n independent observations from an unknown distribution P , a sampling algorithm must output a single observation from a distribution that is close in total variation distance to P while satisfying differential privacy. Sampling abstracts the goal of generating small amounts of realistic-looking data. We provide tight upper and lower bounds for the dataset size needed for this task for three natural families of distributions: arbitrary distributions on $\{1, \dots, k\}$, arbitrary product distributions on $\{0, 1\}^d$, and product distributions on $\{0, 1\}^d$ with bias in each coordinate bounded away from 0 and 1. We demonstrate that, in some parameter regimes, private sampling requires asymptotically fewer observations than learning a description of P nonprivately; in other regimes, however, private sampling proves to be as difficult as private learning. Notably, for some classes of distributions, the overhead in the number of observations needed for private learning compared to non-private learning is completely captured by the number of observations needed for private sampling.

On the Expected Complexity of Maxout Networks

- Hanna Tseran, Guido F. Montufar
- abstract@[open-review](#): Learning with neural networks relies on the complexity of their representable functions, but more importantly, their particular assignment of typical parameters to functions of different complexity. Taking the number of activation regions as a complexity measure, recent works have shown that the practical complexity of deep ReLU networks is often far from the theoretical maximum. In this work, we show that this phenomenon also occurs in networks with maxout (multi-argument) activation functions and when considering the decision boundaries in classification tasks. We also show that the parameter space has a multitude of full-dimensional regions with widely different complexity, and obtain nontrivial lower bounds on the expected complexity. Finally, we investigate different parameter initialization procedures and show that they can increase the speed of convergence in training.

Cross-view Geo-localization with Layer-to-Layer Transformer

- Hongji Yang, Xiufan Lu, Yingying Zhu
- abstract@[open-review](#): In this work, we address the problem of cross-view geo-localization, which estimates the geospatial location of a street view image by matching it with a database of geo-tagged aerial images. The cross-view matching task is extremely challenging due to drastic appearance and geometry differences across views. Unlike existing methods that predominantly fall back on CNN, here we devise a novel layer-to-layer Transformer (L2LTR) that utilizes the properties of self-attention in Transformer to model global dependencies, thus significantly decreasing visual ambiguities in cross-view geo-localization. We also exploit the positional encoding of the Transformer to help the L2LTR understand and correspond geometric configurations between ground and aerial images. Compared to state-of-the-art methods that impose strong assumptions on geometry knowledge, the L2LTR flexibly learns the positional embeddings through the training objective. It hence becomes more practical in many real-world scenarios. Although Transformer is well suited to our task, its vanilla self-attention mechanism independently interacts within image patches in each layer, which overlooks correlations between layers. Instead, this paper proposes a simple yet effective self-cross attention mechanism to improve the quality of learned representations. Self-cross attention models global dependencies between adjacent layers and creates short paths for effective information flow. As a result, the proposed self-cross attention leads to more stable training, improves the generalization ability, and prevents the learned intermediate features from being overly similar. Extensive experiments demonstrate that our L2LTR performs favorably against state-of-the-art methods on standard, fine-grained, and cross-dataset cross-view geo-localization tasks. The code is available online.

TAAC: Temporally Abstract Actor-Critic for Continuous Control

- Haonan Yu, Wei Xu, Haichao Zhang
- abstract@[open-review](#): We present temporally abstract actor-critic (TAAC), a simple but effective off-policy RL algorithm that incorporates closed-loop temporal abstraction into the actor-critic framework. TAAC adds a second-stage binary policy to choose between the previous action and a new action output by an actor. Crucially, its "act-or-repeat" decision hinges on the actually sampled action instead of the expected behavior of the actor. This post-acting switching scheme let the overall policy make more informed decisions. TAAC has two important features: a) persistent exploration, and b) a new compare-through Q operator for multi-step TD backup, specially tailored to the action repetition scenario. We demonstrate TAAC's advantages over several strong baselines across 14 continuous control tasks. Our surprising finding reveals that while achieving top performance, TAAC is able to "mine" a significant number of repeated actions with the trained policy even on continuous tasks whose problem structures on the surface seem to repel action repetition. This suggests that aside from encouraging persistent exploration, action repetition can find its place in a good policy behavior. Code is available at <https://github.com/hnyu/taac>.

Learning Robust Hierarchical Patterns of Human Brain across Many fMRI Studies

- Dushyant Sahoo, Christos Davatzikos
- abstract@[open-review](#): Multi-site fMRI studies face the challenge that the pooling introduces systematic non-biological site-specific variance due to hardware, software, and environment. In this paper, we propose to reduce site-specific variance in the estimation of hierarchical Sparsity Connectivity Patterns (hSCPs) in fMRI data via a simple yet effective matrix factorization while preserving biologically relevant variations. Our method leverages unsupervised adversarial learning to improve the reproducibility of the components. Experiments on simulated datasets display that the proposed method can estimate components with higher accuracy and reproducibility, while preserving age-related variation on a multi-center clinical data set.

Global Convergence to Local Minmax Equilibrium in Classes of Nonconvex Zero-Sum Games

- Tanner Fiez, Lillian Ratliff, Eric Mazumdar, Evan Faulkner, Adhyyan Narang
- abstract@[open-review](#): We study gradient descent-ascent learning dynamics with timescale separation (τ -GDA) in unconstrained continuous action zero-sum games where the minimizing player faces a nonconvex optimization problem and the maximizing player optimizes a Polyak-Lojasiewicz (PL) or strongly-concave (SC) objective. In contrast to past work on gradient-based learning in nonconvex-PL/SC zero-sum games, we assess convergence in relation to natural game-theoretic equilibria instead of only notions of stationarity. In pursuit of this goal, we prove that the only locally stable points of the τ -GDA continuous-time limiting system correspond to strict local minmax equilibria in each class of games. For these classes of games, we exploit timescale separation to construct a potential function that when combined with the stability characterization and an asymptotic saddle avoidance result gives a global asymptotic almost-sure convergence guarantee for the discrete-time gradient descent-ascent update to a set of the strict local minmax equilibrium. Moreover, we provide convergence rates for the gradient descent-ascent dynamics with timescale separation to approximate stationary points.

Bandit Quickest Changepoint Detection

- Aditya Gopalan, Braghadeesh Lakshminarayanan, Venkatesh Saligrama
- abstract@[open-review](#): Many industrial and security applications employ a suite of sensors for detecting abrupt changes in temporal behavior patterns. These abrupt changes typically manifest locally, rendering only a small subset of sensors informative. Continuous monitoring of every sensor can be expensive due to resource constraints, and serves as a motivation for the bandit quickest changepoint detection problem, where sensing actions (or

sensors) are sequentially chosen, and only measurements corresponding to chosen actions are observed. We derive an information-theoretic lower bound on the detection delay for a general class of finitely parameterized probability distributions. We then propose a computationally efficient online sensing scheme, which seamlessly balances the need for exploration of different sensing options with exploitation of querying informative actions. We derive expected delay bounds for the proposed scheme and show that these bounds match our information-theoretic lower bounds at low false alarm rates, establishing optimality of the proposed method. We then perform a number of experiments on synthetic and real datasets demonstrating the effectiveness of our proposed method.

[Can multi-label classification networks know what they don't know?](#)

- Haoran Wang, Weitang Liu, Alex Bocchieri, Yixuan Li
- abstract@[open-review](#): Estimating out-of-distribution (OOD) uncertainty is a major challenge for safely deploying machine learning models in the open-world environment. Improved methods for OOD detection in multi-class classification have emerged, while OOD detection methods for multi-label classification remain underexplored and use rudimentary techniques. We propose JointEnergy, a simple and effective method, which estimates the OOD indicator scores by aggregating label-wise energy scores from multiple labels. We show that JointEnergy can be mathematically interpreted from a joint likelihood perspective. Our results show consistent improvement over previous methods that are based on the maximum-valued scores, which fail to capture joint information from multiple labels. We demonstrate the effectiveness of our method on three common multi-label classification benchmarks, including MS-COCO, PASCAL-VOC, and NUS-WIDE. We show that JointEnergy can reduce the FPR95 by up to 10.05% compared to the previous best baseline, establishing state-of-the-art performance.

[Balanced Chamfer Distance as a Comprehensive Metric for Point Cloud Completion](#)

- Tong Wu, Liang Pan, Junzhe Zhang, Tai WANG, Ziwei Liu, Dahua Lin
- abstract@[open-review](#): Chamfer Distance (CD) and Earth Mover's Distance (EMD) are two broadly adopted metrics for measuring the similarity between two point sets. However, CD is usually insensitive to mismatched local density, and EMD is usually dominated by global distribution while overlooks the fidelity of detailed structures. Besides, their unbounded value range induces a heavy influence from the outliers. These defects prevent them from providing a consistent evaluation. To tackle these problems, we propose a new similarity measure named Density-aware Chamfer Distance (DCD). It is derived from CD and benefits from several desirable properties: 1) it can detect disparity of density distributions and is thus a more intensive measure of similarity compared to CD; 2) it is stricter with detailed structures and significantly more computationally efficient than EMD; 3) the bounded value range encourages a more stable and reasonable evaluation over the whole test set. We adopt DCD to evaluate the point cloud completion task, where experimental results show that DCD pays attention to both the overall structure and local geometric details and provides a more reliable evaluation even when CD and EMD contradict each other. We can also use DCD as the training loss, which outperforms the same model trained with CD loss on all three metrics. In addition, we propose a novel point discriminator module that estimates the priority for another guided down-sampling step, and it achieves noticeable improvements under DCD together with competitive results for both CD and EMD. We hope our work could pave the way for a more comprehensive and practical point cloud similarity evaluation. Our code will be available at https://github.com/wutong16/DensityawareChamfer_Distance.

[Optimal Gradient-based Algorithms for Non-concave Bandit Optimization](#)

- Baihe Huang, Kaixuan Huang, Sham Kakade, Jason D. Lee, Qi Lei, Runzhe Wang, Jiaqi Yang
- abstract@[open-review](#): Bandit problems with linear or concave reward have been extensively studied, but relatively few works have studied bandits with non-concave reward. This work considers a large family of bandit problems where the unknown underlying reward function is non-concave, including the low-rank generalized linear bandit problems and two-layer neural network with polynomial activation bandit problem. For the low-rank generalized linear bandit problem, we provide a minimax-optimal algorithm in the dimension, refuting both conjectures in \cite{lu2021low,jun2019bilinear}. Our algorithms are based on a unified zeroth-order optimization paradigm that applies in great generality and attains optimal rates in several structured polynomial settings (in the dimension). We further demonstrate the applicability of our algorithms in RL in the generative model setting, resulting in improved sample complexity over prior approaches. Finally, we show that the standard optimistic algorithms (e.g., UCB) are sub-optimal by dimension factors. In the neural net setting (with polynomial activation functions) with noiseless reward, we provide a bandit algorithm with sample complexity equal to the intrinsic algebraic dimension. Again, we show that optimistic approaches have worse sample complexity, polynomial in the extrinsic dimension (which could be exponentially worse in the polynomial degree).

[On Optimal Interpolation in Linear Regression](#)

- Eduard Oravkin, Patrick Rebeschini
- abstract@[open-review](#): Understanding when and why interpolating methods generalize well has recently been a topic of interest in statistical learning theory. However, systematically connecting interpolating methods to achievable notions of optimality has only received partial attention. In this paper, we ask the question of what is the optimal way to interpolate in linear regression using functions that are linear in the response variable (as the case for the Bayes optimal estimator in ridge regression) and depend on the data, the population covariance of the data, the signal-to-noise ratio and the covariance of the prior for the signal, but do not depend on the value of the signal itself nor the noise vector in the training data. We provide a closed-form expression for the interpolator that achieves this notion of optimality and show that it can be derived as the limit of preconditioned gradient descent with a specific initialization. We identify a regime where the minimum-norm interpolator provably generalizes arbitrarily worse than the optimal response-linear achievable interpolator that we introduce, and validate with numerical experiments that the notion of optimality we consider can be achieved by interpolating methods that only use the training data as input in the case of an isotropic prior. Finally, we extend the notion of optimal response-linear interpolation to random features regression under a linear data-generating model.

[Differentiable Optimization of Generalized Nondecomposable Functions using Linear Programs](#)

- Zihang Meng, Lopamudra Mukherjee, Yichao Wu, Vikas Singh, Sathya Ravi
- abstract@[open-review](#): We propose a framework which makes it feasible to directly train deep neural networks with respect to popular families of task-specific non-decomposable performance measures such as AUC, multi-class AUC, \$F\$-measure and others. A common feature of the optimization model that emerges from these tasks is that it involves solving a Linear Programs (LP) during training where representations learned by upstream layers characterize the constraints or the feasible set. The constraint matrix is not only large but the constraints are also modified at each iteration. We show how adopting a set of ingenious ideas proposed by Mangasarian for 1-norm SVMs -- which advocates for solving LPs with a generalized Newton method -- provides a simple and effective solution that can be run on the GPU. In particular, this strategy needs little unrolling, which makes it more efficient during backward pass. Further, even when the constraint matrix is too large to fit on the GPU memory (say large minibatch settings), we show that running the Newton method in a lower dimensional space yields accurate gradients for training, by utilizing a statistical concept called {\em sufficient} dimension reduction. While a number of specialized algorithms have been proposed for the models that we describe here, our module turns out to be applicable without any specific adjustments or relaxations. We describe each use case, study its properties and demonstrate the efficacy of the approach over alternatives which use surrogate lower bounds and often, specialized optimization schemes. Frequently, we achieve superior computational behavior and performance improvements on common datasets used in the literature.

[Towards Understanding Cooperative Multi-Agent Q-Learning with Value Factorization](#)

- Jianhao Wang, Zhizhou Ren, Beining Han, Jianing Ye, Chongjie Zhang
- abstract@[open-review](#): Value factorization is a popular and promising approach to scaling up multi-agent reinforcement learning in cooperative settings, which balances the learning scalability and the representational capacity of value functions. However, the theoretical understanding of such methods is limited. In this paper, we formalize a multi-agent fitted Q-iteration framework for analyzing factorized multi-agent Q-learning. Based on this framework, we investigate linear value factorization and reveal that multi-agent Q-learning with this simple decomposition implicitly realizes a powerful counterfactual credit assignment, but may not converge in some settings. Through further analysis, we find that on-policy training or richer joint value function classes can improve its local or global convergence properties, respectively. Finally, to support our theoretical implications in practical realization, we conduct an empirical analysis of state-of-the-art deep multi-agent Q-learning algorithms on didactic examples and a broad set of StarCraft II unit micromanagement tasks.

[**Margin-Independent Online Multiclass Learning via Convex Geometry**](#)

- Guru Guruganesh, Allen Liu, Jon Schneider, Joshua Wang
- abstract@[open-review](#): We consider the problem of multi-class classification, where a stream of adversarially chosen queries arrive and must be assigned a label online. Unlike traditional bounds which seek to minimize the misclassification rate, we minimize the total distance from each query to the region corresponding to its assigned label. When the true labels are determined via a nearest neighbor partition -- i.e. the label of a point is given by which of $\$k$ centers it is closest to in Euclidean distance -- we show that one can achieve a loss that is independent of the total number of queries. We complement this result by showing that learning general convex sets requires an almost linear loss per query. Our results build off of regret guarantees for the problem of contextual search. In addition, we develop a novel reduction technique from multiclass classification to binary classification which may be of independent interest.

[**STEP: Out-of-Distribution Detection in the Presence of Limited In-Distribution Labeled Data**](#)

- Zhi Zhou, Lan-Zhe Guo, Zhanzhan Cheng, Yu-Feng Li, Shiliang Pu
- abstract@[open-review](#): Existing semi-supervised learning (SSL) studies typically assume that unlabeled and test data are drawn from the same distribution as labeled data. However, in many real-world applications, it is desirable to have SSL algorithms that not only classify the samples drawn from the same distribution of labeled data but also detect out-of-distribution (OOD) samples drawn from an unknown distribution. In this paper, we study a setting called semi-supervised OOD detection. Two main challenges compared with previous OOD detection settings are i) the lack of labeled data and in-distribution data; ii) OOD samples could be unseen during training. Efforts on this direction remain limited. In this paper, we present an approach STEP significantly improving OOD detection performance by introducing a new technique: Structure-Keep Unzipping. It learns a new representation space in which OOD samples could be separated well. An efficient optimization algorithm is derived to solve the objective. Comprehensive experiments across various OOD detection benchmarks clearly show that our STEP approach outperforms other methods by a large margin and achieves remarkable detection performance on several benchmarks.

[**Renyi Differential Privacy of The Subsampled Shuffle Model In Distributed Learning**](#)

- Antonios Girgis, Deepesh Data, Suhas Diggavi
- abstract@[open-review](#): We study privacy in a distributed learning framework, where clients collaboratively build a learning model iteratively through interactions with a server from whom we need privacy. Motivated by stochastic optimization and the federated learning (FL) paradigm, we focus on the case where a small fraction of data samples are randomly sub-sampled in each round to participate in the learning process, which also enables privacy amplification. To obtain even stronger local privacy guarantees, we study this in the shuffle privacy model, where each client randomizes its response using a local differentially private (LDP) mechanism and the server only receives a random permutation (shuffle) of the clients' responses without their association to each client. The principal result of this paper is a privacy-optimization performance trade-off for discrete randomization mechanisms in this sub-sampled shuffle privacy model. This is enabled through a new theoretical technique to analyze the Renyi Differential Privacy (RDP) of the sub-sampled shuffle model. We numerically demonstrate that, for important regimes, with composition our bound yields significant improvement in privacy guarantee over the state-of-the-art approximate Differential Privacy (DP) guarantee (with strong composition) for sub-sampled shuffled models. We also demonstrate numerically significant improvement in privacy-learning performance operating point using real data sets. Despite these advances, an open question is to bridge the gap between lower and upper privacy bounds in our RDP analysis.

[**Gradient-based Editing of Memory Examples for Online Task-free Continual Learning**](#)

- Xisen Jin, Arka Sadhu, Junyi Du, Xiang Ren
- abstract@[open-review](#): We explore task-free continual learning (CL), in which a model is trained to avoid catastrophic forgetting in the absence of explicit task boundaries or identities. Among many efforts on task-free CL, a notable family of approaches are memory-based that store and replay a subset of training examples. However, the utility of stored seen examples may diminish over time since CL models are continually updated. Here, we propose Gradient based Memory EDiting (GMED), a framework for editing stored examples in continuous input space via gradient updates, in order to create more "challenging" examples for replay. GMED-edited examples remain similar to their unedited forms, but can yield increased loss in the upcoming model updates, thereby making the future replays more effective in overcoming catastrophic forgetting. By construction, GMED can be seamlessly applied in conjunction with other memory-based CL algorithms to bring further improvement. Experiments validate the effectiveness of GMED, and our best method significantly outperforms baselines and previous state-of-the-art on five out of six datasets.

[**Tailoring: encoding inductive biases by optimizing unsupervised objectives at prediction time**](#)

- Ferran Alet, Maria Bauza, Kenji Kawaguchi, Nurullah Giray Kuru, Tomás Lozano-Pérez, Leslie Kaelbling
- abstract@[open-review](#): From CNNs to attention mechanisms, encoding inductive biases into neural networks has been a fruitful source of improvement in machine learning. Adding auxiliary losses to the main objective function is a general way of encoding biases that can help networks learn better representations. However, since auxiliary losses are minimized only on training data, they suffer from the same generalization gap as regular task losses. Moreover, by adding a term to the loss function, the model optimizes a different objective than the one we care about. In this work we address both problems: first, we take inspiration from transductive learning and note that after receiving an input but before making a prediction, we can fine-tune our networks on any unsupervised loss. We call this process tailoring, because we customize the model to each input to ensure our prediction satisfies the inductive bias. Second, we formulate meta-tailoring, a nested optimization similar to that in meta-learning, and train our models to perform well on the task objective after adapting them using an unsupervised loss. The advantages of tailoring and meta-tailoring are discussed theoretically and demonstrated empirically on a diverse set of examples.

[**Implicit Bias of SGD for Diagonal Linear Networks: a Provable Benefit of Stochasticity**](#)

- Scott Pesme, Lucas Pillaud-Vivien, Nicolas Flammarion
- abstract@[open-review](#): Understanding the implicit bias of training algorithms is of crucial importance in order to explain the success of overparametrised neural networks. In this paper, we study the dynamics of stochastic gradient descent over diagonal linear networks through its continuous time version, namely stochastic gradient flow. We explicitly characterise the solution chosen by the stochastic flow and prove that it always enjoys better generalisation properties than that of gradient flow. Quite surprisingly, we show that the convergence speed of the training loss controls the magnitude of the biasing effect: the slower the convergence, the better the bias. To fully complete our analysis, we provide convergence guarantees for the dynamics. We also give

experimental results which support our theoretical claims. Our findings highlight the fact that structured noise can induce better generalisation and they help explain the greater performances of stochastic gradient descent over gradient descent observed in practice.

[Iterative Teacher-Aware Learning](#)

- Luyao Yuan, Dongruo Zhou, Junhong Shen, Jingdong Gao, Jeffrey L Chen, Quanquan Gu, Ying Nian Wu, Song-Chun Zhu
- abstract@[open-review](#): In human pedagogy, teachers and students can interact adaptively to maximize communication efficiency. The teacher adjusts her teaching method for different students, and the student, after getting familiar with the teacher's instruction mechanism, can infer the teacher's intention to learn faster. Recently, the benefits of integrating this cooperative pedagogy into machine concept learning in discrete spaces have been proved by multiple works. However, how cooperative pedagogy can facilitate machine parameter learning hasn't been thoroughly studied. In this paper, we propose a gradient optimization based teacher-aware learner who can incorporate teacher's cooperative intention into the likelihood function and learn provably faster compared with the naive learning algorithms used in previous machine teaching works. We give theoretical proof that the iterative teacher-aware learning (ITAL) process leads to local and global improvements. We then validate our algorithms with extensive experiments on various tasks including regression, classification, and inverse reinforcement learning using synthetic and real data. We also show the advantage of modeling teacher-awareness when agents are learning from human teachers.

[Clockwork Variational Autoencoders](#)

- Vaibhav Saxena, Jimmy Ba, Danijar Hafner
- abstract@[open-review](#): Deep learning has enabled algorithms to generate realistic images. However, accurately predicting long video sequences requires understanding long-term dependencies and remains an open challenge. While existing video prediction models succeed at generating sharp images, they tend to fail at accurately predicting far into the future. We introduce the Clockwork VAE (CW-VAE), a video prediction model that leverages a hierarchy of latent sequences, where higher levels tick at slower intervals. We demonstrate the benefits of both hierarchical latents and temporal abstraction on 4 diverse video prediction datasets with sequences of up to 1000 frames, where CW-VAE outperforms top video prediction models. Additionally, we propose a Minecraft benchmark for long-term video prediction. We conduct several experiments to gain insights into CW-VAE and confirm that slower levels learn to represent objects that change more slowly in the video, and faster levels learn to represent faster objects.

[How Does it Sound?](#)

- Kun Su, Xiulong Liu, Eli Shlizerman
- abstract@[open-review](#): One of the primary purposes of video is to capture people and their unique activities. It is often the case that the experience of watching the video can be enhanced by adding a musical soundtrack that is in-sync with the rhythmic features of these activities. How would this soundtrack sound? Such a problem is challenging since little is known about capturing the rhythmic nature of free body movements. In this work, we explore this problem and propose a novel system, called 'RhythmicNet', which takes as an input a video which includes human movements and generates a soundtrack for it. RhythmicNet works directly with human movements by extracting skeleton keypoints and implements a sequence of models which translate the keypoints to rhythmic sounds. RhythmicNet follows the natural process of music improvisation which includes the prescription of streams of the beat, the rhythm and the melody. In particular, RhythmicNet first infers the music beat and the style pattern from body keypoints per each frame to produce rhythm. Next, it implements a transformer-based model to generate the hits of drum instruments and implements a U-net based model to generate the velocity and the offsets of the instruments. Additional types of instruments are added to the soundtrack by further conditioning on the generated drum sounds. We evaluate RhythmicNet on large scale datasets of videos that include body movements with inherit sound association, such as dance, as well as 'in the wild' internet videos of various movements and actions. We show that the method can generate plausible music that aligns well with different types of human movements.

[Stabilizing Dynamical Systems via Policy Gradient Methods](#)

- Juan Perdomo, Jack Umenberger, Max Simchowitz
- abstract@[open-review](#): Stabilizing an unknown control system is one of the most fundamental problems in control systems engineering. In this paper, we provide a simple, model-free algorithm for stabilizing fully observed dynamical systems. While model-free methods have become increasingly popular in practice due to their simplicity and flexibility, stabilization via direct policy search has received surprisingly little attention. Our algorithm proceeds by solving a series of discounted LQR problems, where the discount factor is gradually increased. We prove that this method efficiently recovers a stabilizing controller for linear systems, and for smooth, nonlinear systems within a neighborhood of their equilibria. Our approach overcomes a significant limitation of prior work, namely the need for a pre-given stabilizing control policy. We empirically evaluate the effectiveness of our approach on common control benchmarks.

[Language models enable zero-shot prediction of the effects of mutations on protein function](#)

- Joshua Meier, Roshan Rao, Robert Verkuil, Jason Liu, Tom Sercu, Alex Rives
- abstract@[open-review](#): Modeling the effect of sequence variation on function is a fundamental problem for understanding and designing proteins. Since evolution encodes information about function into patterns in protein sequences, unsupervised models of variant effects can be learned from sequence data. The approach to date has been to fit a model to a family of related sequences. The conventional setting is limited, since a new model must be trained for each prediction task. We show that using only zero-shot inference, without any supervision from experimental data or additional training, protein language models capture the functional effects of sequence variation, performing at state-of-the-art.

[Deep Reinforcement Learning at the Edge of the Statistical Precipice](#)

- Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C. Courville, Marc Bellemare
- abstract@[open-review](#): Deep reinforcement learning (RL) algorithms are predominantly evaluated by comparing their relative performance on a large suite of tasks. Most published results on deep RL benchmarks compare point estimates of aggregate performance such as mean and median scores across tasks, ignoring the statistical uncertainty implied by the use of a finite number of training runs. Beginning with the Arcade Learning Environment (ALE), the shift towards computationally-demanding benchmarks has led to the practice of evaluating only a small number of runs per task, exacerbating the statistical uncertainty in point estimates. In this paper, we argue that reliable evaluation in the few run deep RL regime cannot ignore the uncertainty in results without running the risk of slowing down progress in the field. We illustrate this point using a case study on the Atari 100k benchmark, where we find substantial discrepancies between conclusions drawn from point estimates alone versus a more thorough statistical analysis. With the aim of increasing the field's confidence in reported results with a handful of runs, we advocate for reporting interval estimates of aggregate performance and propose performance profiles to account for the variability in results, as well as present more robust and efficient aggregate metrics, such as interquartile mean scores, to achieve small uncertainty in results. Using such statistical tools, we scrutinize performance evaluations of existing algorithms on other widely used RL benchmarks including the ALE, Procgen, and the DeepMind Control Suite, again revealing discrepancies in prior comparisons. Our findings call for a change in how we evaluate performance in deep RL, for which we present a more rigorous evaluation methodology, accompanied with an open-source library rliable, to prevent unreliable results from stagnating the field. This work received an outstanding paper award at NeurIPS 2021.

[DRONE: Data-aware Low-rank Compression for Large NLP Models](#)

- Patrick Chen, Hsiang-Fu Yu, Inderjit Dhillon, Cho-Jui Hsieh
- abstract@[open-review](#): The representations learned by large-scale NLP models such as BERT have been widely used in various tasks. However, the increasing model size of the pre-trained models also brings efficiency challenges, including inference speed and model size when deploying models on mobile devices. Specifically, most operations in BERT consist of matrix multiplications. These matrices are not low-rank and thus canonical matrix decomposition could not find an efficient approximation. In this paper, we observe that the learned representation of each layer lies in a low-dimensional space. Based on this observation, we propose DRONE (data-aware low-rank compression), a provably optimal low-rank decomposition of weight matrices, which has a simple closed form solution that can be efficiently computed. DRONE can be applied to both fully connected and self-attention layers appearing in the BERT model. In addition to compressing standard models, our method can also be used on distilled BERT models to further improve compression rate. Experimental results show that DRONE is able to improve both model size and inference speed with limited loss in accuracy. Specifically, DRONE alone achieves 1.92x speedup on the MRPC task with only 1.5% loss in accuracy, and when DRONE is combined with distillation, it further achieves over 12.3x speedup on various natural language inference tasks.

[DSelect-k: Differentiable Selection in the Mixture of Experts with Applications to Multi-Task Learning](#)

- Hussein Hazimeh, Zhe Zhao, Aakanksha Chowdhery, Maheswaran Sathiamoorthy, Yihua Chen, Rahul Mazumder, Lichan Hong, Ed Chi
- abstract@[open-review](#): The Mixture-of-Experts (MoE) architecture is showing promising results in improving parameter sharing in multi-task learning (MTL) and in scaling high-capacity neural networks. State-of-the-art MoE models use a trainable "sparse gate" to select a subset of the experts for each input example. While conceptually appealing, existing sparse gates, such as Top-k, are not smooth. The lack of smoothness can lead to convergence and statistical performance issues when training with gradient-based methods. In this paper, we develop DSelect-k: a continuously differentiable and sparse gate for MoE, based on a novel binary encoding formulation. The gate can be trained using first-order methods, such as stochastic gradient descent, and offers explicit control over the number of experts to select. We demonstrate the effectiveness of DSelect-k on both synthetic and real MTL datasets with up to 128 tasks. Our experiments indicate that DSelect-k can achieve statistically significant improvements in prediction and expert selection over popular MoE gates. Notably, on a real-world, large-scale recommender system, DSelect-k achieves over 22% improvement in predictive performance compared to Top-k. We provide an open-source implementation of DSelect-k.

[Mind the Gap: Assessing Temporal Generalization in Neural Language Models](#)

- Angeliki Lazaridou, Adhi Kuncoro, Elena Gribovskaya, Devang Agrawal, Adam Liska, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d'Autume, Tomas Kociský, Sebastian Ruder, Dani Yogatama, Kris Cao, Susannah Young, Phil Blunsom
- abstract@[open-review](#): Our world is open-ended, non-stationary, and constantly evolving; thus what we talk about and how we talk about it change over time. This inherent dynamic nature of language contrasts with the current static language modelling paradigm, which trains and evaluates models on utterances from overlapping time periods. Despite impressive recent progress, we demonstrate that Transformer-XL language models perform worse in the realistic setup of predicting future utterances from beyond their training period, and that model performance becomes increasingly worse with time. We find that, while increasing model size alone—a key driver behind recent progress—does not solve this problem, having models that continually update their knowledge with new information can indeed mitigate this performance degradation over time. Hence, given the compilation of ever-larger language modelling datasets, combined with the growing list of language-model-based NLP applications that require up-to-date factual knowledge about the world, we argue that now is the right time to rethink the static way in which we currently train and evaluate our language models, and develop adaptive language models that can remain up-to-date with respect to our ever-changing and non-stationary world. We publicly release our dynamic, streaming language modelling benchmarks for WMT and arXiv to facilitate language model evaluation that takes temporal dynamics into account.

[Heavy Tails in SGD and Compressibility of Overparametrized Neural Networks](#)

- Melih Barsbey, Milad Sefidgaran, Murat A. Erdogdu, Gaël Richard, Umut Simsekli
- abstract@[open-review](#): Neural network compression techniques have become increasingly popular as they can drastically reduce the storage and computation requirements for very large networks. Recent empirical studies have illustrated that even simple pruning strategies can be surprisingly effective, and several theoretical studies have shown that compressible networks (in specific senses) should achieve a low generalization error. Yet, a theoretical characterization of the underlying causes that make the networks amenable to such simple compression schemes is still missing. In this study, focusing our attention on stochastic gradient descent (SGD), our main contribution is to link compressibility to two recently established properties of SGD: (i) as the network size goes to infinity, the system can converge to a mean-field limit, where the network weights behave independently [DBDFŞ20], (ii) for a large step-size/batch-size ratio, the SGD iterates can converge to a heavy-tailed stationary distribution [HM20, GSZ21]. Assuming that both of these phenomena occur simultaneously, we prove that the networks are guaranteed to be '\$\ell_p\$-compressible', and the compression errors of different pruning techniques (magnitude, singular value, or node pruning) become arbitrarily small as the network size increases. We further prove generalization bounds adapted to our theoretical framework, which are consistent with the observation that the generalization error will be lower for more compressible networks. Our theory and numerical study on various neural networks show that large step-size/batch-size ratios introduce heavy tails, which, in combination with overparametrization, result in compressibility.

[Targeted Neural Dynamical Modeling](#)

- Cole Hurwitz, Akash Srivastava, Kai Xu, Justin Jude, Matthew Perich, Lee Miller, Matthias Hennig
- abstract@[open-review](#): Latent dynamics models have emerged as powerful tools for modeling and interpreting neural population activity. Recently, there has been a focus on incorporating simultaneously measured behaviour into these models to further disentangle sources of neural variability in their latent space. These approaches, however, are limited in their ability to capture the underlying neural dynamics (e.g. linear) and in their ability to relate the learned dynamics back to the observed behaviour (e.g. no time lag). To this end, we introduce Targeted Neural Dynamical Modeling (TNDM), a nonlinear state-space model that jointly models the neural activity and external behavioural variables. TNDM decomposes neural dynamics into behaviourally relevant and behaviourally irrelevant dynamics; the relevant dynamics are used to reconstruct the behaviour through a flexible linear decoder and both sets of dynamics are used to reconstruct the neural activity through a linear decoder with no time lag. We implement TNDM as a sequential variational autoencoder and validate it on simulated recordings and recordings taken from the premotor and motor cortex of a monkey performing a center-out reaching task. We show that TNDM is able to learn low-dimensional latent dynamics that are highly predictive of behaviour without sacrificing its fit to the neural data.

[Exploiting the Intrinsic Neighborhood Structure for Source-free Domain Adaptation](#)

- Shiqi Yang, Yaxing Wang, Joost van de Weijer, Luis Herranz, Shangling Jui
- abstract@[open-review](#): Domain adaptation (DA) aims to alleviate the domain shift between source domain and target domain. Most DA methods require access to the source data, but often that is not possible (e.g. due to data privacy or intellectual property). In this paper, we address the challenging source-free domain adaptation (SFDA) problem, where the source pretrained model is adapted to the target domain in the absence of source data. Our method is based on the observation that target data, which might no longer align with the source domain classifier, still forms clear clusters. We capture this intrinsic structure by defining local affinity of the target data, and encourage label consistency among data with high local affinity. We observe that higher affinity should be assigned to reciprocal neighbors, and propose a self regularization loss to decrease the negative impact of noisy neighbors. Furthermore, to aggregate information with more context, we consider expanded neighborhoods with small affinity values. In the experimental results we verify that the inherent structure of the target features is an important source of information for domain adaptation. We demonstrate that this local structure can be

efficiently captured by considering the local neighbors, the reciprocal neighbors, and the expanded neighborhood. Finally, we achieve state-of-the-art performance on several 2D image and 3D point cloud recognition datasets. Code is available in https://github.com/Albert0147/SFDA_neighbors.

[Learning with Noisy Correspondence for Cross-modal Matching](#)

- Zhenyu Huang, Guocheng Niu, Xiao Liu, Wenbiao Ding, Xinyan Xiao, Hua Wu, Xi Peng
- abstract@[open-review](#): Cross-modal matching, which aims to establish the correspondence between two different modalities, is fundamental to a variety of tasks such as cross-modal retrieval and vision-and-language understanding. Although a huge number of cross-modal matching methods have been proposed and achieved remarkable progress in recent years, almost all of these methods implicitly assume that the multimodal training data are correctly aligned. In practice, however, such an assumption is extremely expensive even impossible to satisfy. Based on this observation, we reveal and study a latent and challenging direction in cross-modal matching, named noisy correspondence, which could be regarded as a new paradigm of noisy labels. Different from the traditional noisy labels which mainly refer to the errors in category labels, our noisy correspondence refers to the mismatch paired samples. To solve this new problem, we propose a novel method for learning with noisy correspondence, named Noisy Correspondence Rectifier (NCR). In brief, NCR divides the data into clean and noisy partitions based on the memorization effect of neural networks and then rectifies the correspondence via an adaptive prediction model in a co-teaching manner. To verify the effectiveness of our method, we conduct experiments by using the image-text matching as a showcase. Extensive experiments on Flickr30K, MS-COCO, and Conceptual Captions verify the effectiveness of our method. The code could be accessed from www.pengxi.me .

[Offline Reinforcement Learning with Reverse Model-based Imagination](#)

- Jianhao Wang, Wenzhe Li, Haozhe Jiang, Guangxiang Zhu, Siyuan Li, Chongjie Zhang
- abstract@[open-review](#): In offline reinforcement learning (offline RL), one of the main challenges is to deal with the distributional shift between the learning policy and the given dataset. To address this problem, recent offline RL methods attempt to introduce conservatism bias to encourage learning in high-confidence areas. Model-free approaches directly encode such bias into policy or value function learning using conservative regularizations or special network structures, but their constrained policy search limits the generalization beyond the offline dataset. Model-based approaches learn forward dynamics models with conservatism quantifications and then generate imaginary trajectories to extend the offline datasets. However, due to limited samples in offline datasets, conservatism quantifications often suffer from overgeneralization in out-of-support regions. The unreliable conservative measures will mislead forward model-based imaginations to undesired areas, leading to overaggressive behaviors. To encourage more conservatism, we propose a novel model-based offline RL framework, called Reverse Offline Model-based Imagination (ROMI). We learn a reverse dynamics model in conjunction with a novel reverse policy, which can generate rollouts leading to the target goal states within the offline dataset. These reverse imaginations provide informed data augmentation for model-free policy learning and enable conservative generalization beyond the offline dataset. ROMI can effectively combine with off-the-shelf model-free algorithms to enable model-based generalization with proper conservatism. Empirical results show that our method can generate more conservative behaviors and achieve state-of-the-art performance on offline RL benchmark tasks.

[Parameter Prediction for Unseen Deep Architectures](#)

- Boris Knyazev, Michal Drozdzał, Graham W. Taylor, Adriana Romero Soriano
- abstract@[open-review](#): Deep learning has been successful in automating the design of features in machine learning pipelines. However, the algorithms optimizing neural network parameters remain largely hand-designed and computationally inefficient. We study if we can use deep learning to directly predict these parameters by exploiting the past knowledge of training other networks. We introduce a large-scale dataset of diverse computational graphs of neural architectures - DeepNets-1M - and use it to explore parameter prediction on CIFAR-10 and ImageNet. By leveraging advances in graph neural networks, we propose a hypernetwork that can predict performant parameters in a single forward pass taking a fraction of a second, even on a CPU. The proposed model achieves surprisingly good performance on unseen and diverse networks. For example, it is able to predict all 24 million parameters of a ResNet-50 achieving a 60% accuracy on CIFAR-10. On ImageNet, top-5 accuracy of some of our networks approaches 50%. Our task along with the model and results can potentially lead to a new, more computationally efficient paradigm of training networks. Our model also learns a strong representation of neural architectures enabling their analysis.

[FMMformer: Efficient and Flexible Transformer via Decomposed Near-field and Far-field Attention](#)

- Tan Nguyen, Vai Suliafu, Stanley Osher, Long Chen, Bao Wang
- abstract@[open-review](#): We propose FMMformers, a class of efficient and flexible transformers inspired by the celebrated fast multipole method (FMM) for accelerating interacting particle simulation. FMM decomposes particle-particle interaction into near-field and far-field components and then performs direct and coarse-grained computation, respectively. Similarly, FMMformers decompose the attention into near-field and far-field attention, modeling the near-field attention by a banded matrix and the far-field attention by a low-rank matrix. Computing the attention matrix for FMMformers requires linear complexity in computational time and memory footprint with respect to the sequence length. In contrast, standard transformers suffer from quadratic complexity. We analyze and validate the advantage of FMMformers over the standard transformer on the Long Range Arena and language modeling benchmarks. FMMformers can even outperform the standard transformer in terms of accuracy by a significant margin. For instance, FMMformers achieve an average classification accuracy of \$60.74\%\$ over the five Long Range Arena tasks, which is significantly better than the standard transformer's average accuracy of \$58.70\%\$.

[Square Root Principal Component Pursuit: Tuning-Free Noisy Robust Matrix Recovery](#)

- Junhui Zhang, Jingkai Yan, John Wright
- abstract@[open-review](#): We propose a new framework -- Square Root Principal Component Pursuit -- for low-rank matrix recovery from observations corrupted with noise and outliers. Inspired by the square root Lasso, this new formulation does not require prior knowledge of the noise level. We show that a single, universal choice of the regularization parameter suffices to achieve reconstruction error proportional to the (a priori unknown) noise level. In comparison, previous formulations such as stable PCP rely on noise-dependent parameters to achieve similar performance, and are therefore challenging to deploy in applications where the noise level is unknown. We validate the effectiveness of our new method through experiments on simulated and real datasets. Our simulations corroborate the claim that a universal choice of the regularization parameter yields near optimal performance across a range of noise levels, indicating that the proposed method outperforms the (somewhat loose) bound proved here.

[Neural Bellman-Ford Networks: A General Graph Neural Network Framework for Link Prediction](#)

- Zhaocheng Zhu, Zuobai Zhang, Louis-Pascal Xhonneux, Jian Tang
- abstract@[open-review](#): Link prediction is a very fundamental task on graphs. Inspired by traditional path-based methods, in this paper we propose a general and flexible representation learning framework based on paths for link prediction. Specifically, we define the representation of a pair of nodes as the generalized sum of all path representations, with each path representation as the generalized product of the edge representations in the path. Motivated by the Bellman-Ford algorithm for solving the shortest path problem, we show that the proposed path formulation can be efficiently solved by the generalized Bellman-Ford algorithm. To further improve the capacity of the path formulation, we propose the Neural Bellman-Ford Network (NBFNet), a general graph neural network framework that solves the path formulation with learned operators in the generalized Bellman-Ford algorithm. The NBFNet parameterizes the generalized Bellman-Ford algorithm with 3 neural components, namely Indicator, Message and Aggregate functions, which corresponds to the boundary condition, multiplication operator, and summation operator respectively. The NBFNet covers many traditional path-based methods, and

can be applied to both homogeneous graphs and multi-relational graphs (e.g., knowledge graphs) in both transductive and inductive settings. Experiments on both homogeneous graphs and knowledge graphs show that the proposed NBFNet outperforms existing methods by a large margin in both transductive and inductive settings, achieving new state-of-the-art results.

[CorticalFlow: A Diffeomorphic Mesh Transformer Network for Cortical Surface Reconstruction](#)

- Leo Lebrat, Rodrigo Santa Cruz, Frederic de Gournay, Darren Fu, Pierrick Bourgeat, Jurgen Fripp, Clinton Fookes, Olivier Salvado
- abstract@[open-review](#): In this paper, we introduce CorticalFlow, a new geometric deep-learning model that, given a 3-dimensional image, learns to deform a reference template towards a targeted object. To conserve the template mesh's topological properties, we train our model over a set of diffeomorphic transformations. This new implementation of a flow Ordinary Differential Equation (ODE) framework benefits from a small GPU memory footprint, allowing the generation of surfaces with several hundred thousand vertices. To reduce topological errors introduced by its discrete resolution, we derive numeric conditions which improve the manifoldness of the predicted triangle mesh. To exhibit the utility of CorticalFlow, we demonstrate its performance for the challenging task of brain cortical surface reconstruction. In contrast to the current state-of-the-art, CorticalFlow produces superior surfaces while reducing the computation time from nine and a half minutes to one second. More significantly, CorticalFlow enforces the generation of anatomically plausible surfaces; the absence of which has been a major impediment restricting the clinical relevance of such surface reconstruction methods.

[Bridging the Gap Between Practice and PAC-Bayes Theory in Few-Shot Meta-Learning](#)

- Nan Ding, Xi Chen, Tomer Levinboim, Sebastian Goodman, Radu Soricut
- abstract@[open-review](#): Despite recent advances in its theoretical understanding, there still remains a significant gap in the ability of existing PAC-Bayesian theories on meta-learning to explain performance improvements in the few-shot learning setting, where the number of training examples in the target tasks is severely limited. This gap originates from an assumption in the existing theories which supposes that the number of training examples in the observed tasks and the number of training examples in the target tasks follow the same distribution, an assumption that rarely holds in practice. By relaxing this assumption, we develop two PAC-Bayesian bounds tailored for the few-shot learning setting and show that two existing meta-learning algorithms (MAML and Reptile) can be derived from our bounds, thereby bridging the gap between practice and PAC-Bayesian theories. Furthermore, we derive a new computationally-efficient PACMAML algorithm, and show it outperforms existing meta-learning algorithms on several few-shot benchmark datasets.

[SLOE: A Faster Method for Statistical Inference in High-Dimensional Logistic Regression](#)

- Steve Yadlowsky, Taedong Yun, Cory Y McLean, Alexander D'Amour
- abstract@[open-review](#): Logistic regression remains one of the most widely used tools in applied statistics, machine learning and data science. However, in moderately high-dimensional problems, where the number of features d is a non-negligible fraction of the sample size n , the logistic regression maximum likelihood estimator (MLE), and statistical procedures based the large-sample approximation of its distribution, behave poorly. Recently, Sur and Candès (2019) showed that these issues can be corrected by applying a new approximation of the MLE's sampling distribution in this high-dimensional regime. Unfortunately, these corrections are difficult to implement in practice, because they require an estimate of the signal strength, which is a function of the underlying parameters β of the logistic regression. To address this issue, we propose SLOE, a fast and straightforward approach to estimate the signal strength in logistic regression. The key insight of SLOE is that the Sur and Candès (2019) correction can be reparameterized in terms of the corrupted signal strength, which is only a function of the estimated parameters $\widehat{\beta}$. We propose an estimator for this quantity, prove that it is consistent in the relevant high-dimensional regime, and show that dimensionality correction using SLOE is accurate in finite samples. Compared to the existing ProbeFrontier heuristic, SLOE is conceptually simpler and orders of magnitude faster, making it suitable for routine use. We demonstrate the importance of routine dimensionality correction in the Heart Disease dataset from the UCI repository, and a genomics application using data from the UK Biobank.

[ELLA: Exploration through Learned Language Abstraction](#)

- Suvir Mirchandani, Siddharth Karamcheti, Dorsa Sadigh
- abstract@[open-review](#): Building agents capable of understanding language instructions is critical to effective and robust human-AI collaboration. Recent work focuses on training these agents via reinforcement learning in environments with synthetic language; however, instructions often define long-horizon, sparse-reward tasks, and learning policies requires many episodes of experience. We introduce ELLA: Exploration through Learned Language Abstraction, a reward shaping approach geared towards boosting sample efficiency in sparse reward environments by correlating high-level instructions with simpler low-level constituents. ELLA has two key elements: 1) A termination classifier that identifies when agents complete low-level instructions, and 2) A relevance classifier that correlates low-level instructions with success on high-level tasks. We learn the termination classifier offline from pairs of instructions and terminal states. Notably, in departure from prior work in language and abstraction, we learn the relevance classifier online, without relying on an explicit decomposition of high-level instructions to low-level instructions. On a suite of complex BabyAI environments with varying instruction complexities and reward sparsity, ELLA shows gains in sample efficiency relative to language-based shaping and traditional RL methods.

[Learning Distilled Collaboration Graph for Multi-Agent Perception](#)

- Yiming Li, Shunli Ren, Pengxiang Wu, Siheng Chen, Chen Feng, Wenjun Zhang
- abstract@[open-review](#): To promote better performance-bandwidth trade-off for multi-agent perception, we propose a novel distilled collaboration graph (DiscoGraph) to model trainable, pose-aware, and adaptive collaboration among agents. Our key novelties lie in two aspects. First, we propose a teacher-student framework to train DiscoGraph via knowledge distillation. The teacher model employs an early collaboration with holistic-view inputs; the student model is based on intermediate collaboration with single-view inputs. Our framework trains DiscoGraph by constraining post-collaboration feature maps in the student model to match the correspondences in the teacher model. Second, we propose a matrix-valued edge weight in DiscoGraph. In such a matrix, each element reflects the inter-agent attention at a specific spatial region, allowing an agent to adaptively highlight the informative regions. During inference, we only need to use the student model named as the distilled collaboration network (DiscoNet). Attributed to the teacher-student framework, multiple agents with the shared DiscoNet could collaboratively approach the performance of a hypothetical teacher model with a holistic view. Our approach is validated on V2X-Sim 1.0, a large-scale multi-agent perception dataset that we synthesized using CARLA and SUMO co-simulation. Our quantitative and qualitative experiments in multi-agent 3D object detection show that DiscoNet could not only achieve a better performance-bandwidth trade-off than the state-of-the-art collaborative perception methods, but also bring more straightforward design rationale. Our code is available on <https://github.com/ai4ce/DiscoNet>.

[Federated-EM with heterogeneity mitigation and variance reduction](#)

- Aymeric Dieuleveut, Gersende Fort, Eric Moulines, Geneviève Robin
- abstract@[open-review](#): The Expectation Maximization (EM) algorithm is the default algorithm for inference in latent variable models. As in any other field of machine learning, applications of latent variable models to very large datasets make the use of advanced parallel and distributed architecture mandatory. This paper introduces FedEM, which is the first extension of the EM algorithm to the federated learning context. FedEM is a new communication efficient method, which handles partial participation of local devices, and is robust to heterogeneous distribution of the datasets. To alleviate the communication bottleneck, FedEM compresses appropriately defined complete data sufficient statistics. We also develop and analyze an

extension of FedEM to further incorporate a variance reduction scheme. In all cases, we derive finite-time complexity bounds for smooth non-convex problems. Numerical results are presented to support our theoretical findings, as well as an application to federated missing values imputation for biodiversity monitoring.

[On the Role of Optimization in Double Descent: A Least Squares Study](#)

- Ilja Kuzborskij, Csaba Szepesvari, Omar Rivasplata, Amal Rannen-Triki, Razvan Pascanu
- abstract@[open-review](#): Empirically it has been observed that the performance of deep neural networks steadily improves with increased model size, contradicting the classical view on overfitting and generalization. Recently, the double descent phenomenon has been proposed to reconcile this observation with theory, suggesting that the test error has a second descent when the model becomes sufficiently overparameterized, as the model size itself acts as an implicit regularizer. In this paper we add to the growing body of work in this space, providing a careful study of learning dynamics as a function of model size for the least squares scenario. We show an excess risk bound for the gradient descent solution of the least squares objective. The bound depends on the smallest non-zero eigenvalue of the sample covariance matrix of the input features, via a functional form that has the double descent behaviour. This gives a new perspective on the double descent curves reported in the literature, as our analysis of the excess risk allows to decouple the effect of optimization and generalization error. In particular, we find that in the case of noiseless regression, double descent is explained solely by optimization-related quantities, which was missed in studies focusing on the Moore-Penrose pseudoinverse solution. We believe that our derivation provides an alternative view compared to existing works, shedding some light on a possible cause of this phenomenon, at least in the considered least squares setting. We empirically explore if our predictions hold for neural networks, in particular whether the spectrum of the sample covariance of features at intermediary hidden layers has a similar behaviour as the one predicted by our derivations in the least squares setting.

[Neural Architecture Dilution for Adversarial Robustness](#)

- Yanxi Li, Zhaohui Yang, Yunhe Wang, Chang Xu
- abstract@[open-review](#): With the tremendous advances in the architecture and scale of convolutional neural networks (CNNs) over the past few decades, they can easily reach or even exceed the performance of humans in certain tasks. However, a recently discovered shortcoming of CNNs is that they are vulnerable to adversarial attacks. Although the adversarial robustness of CNNs can be improved by adversarial training, there is a trade-off between standard accuracy and adversarial robustness. From the neural architecture perspective, this paper aims to improve the adversarial robustness of the backbone CNNs that have a satisfactory accuracy. Under a minimal computational overhead, the introduction of a dilation architecture is expected to be friendly with the standard performance of the backbone CNN while pursuing adversarial robustness. Theoretical analyses on the standard and adversarial error bounds naturally motivate the proposed neural architecture dilation algorithm. Experimental results on real-world datasets and benchmark neural networks demonstrate the effectiveness of the proposed algorithm to balance the accuracy and adversarial robustness.

[Clustering Effect of Adversarial Robust Models](#)

- Yang Bai, Xin Yan, Yong Jiang, Shu-Tao Xia, Yisen Wang
- abstract@[open-review](#): Adversarial robustness has received increasing attention along with the study of adversarial examples. So far, existing works show that robust models not only obtain robustness against various adversarial attacks but also boost the performance in some downstream tasks. However, the underlying mechanism of adversarial robustness is still not clear. In this paper, we interpret adversarial robustness from the perspective of linear components, and find that there exist some statistical properties for comprehensively robust models. Specifically, robust models show obvious hierarchical clustering effect on their linearized sub-networks, when removing or replacing all non-linear components (e.g., batch normalization, maximum pooling, or activation layers). Based on these observations, we propose a novel understanding of adversarial robustness and apply it on more tasks including domain adaption and robustness boosting. Experimental evaluations demonstrate the rationality and superiority of our proposed clustering strategy. Our code is available at <https://github.com/bymavis/AdvWeightNeurIPS2021>.

[On the Cryptographic Hardness of Learning Single Periodic Neurons](#)

- Min Jae Song, Ilias Zadik, Joan Bruna
- abstract@[open-review](#): We show a simple reduction which demonstrates the cryptographic hardness of learning a single periodic neuron over isotropic Gaussian distributions in the presence of noise. More precisely, our reduction shows that any polynomial-time algorithm (not necessarily gradient-based) for learning such functions under small noise implies a polynomial-time quantum algorithm for solving worst-case lattice problems, whose hardness form the foundation of lattice-based cryptography. Our core hard family of functions, which are well-approximated by one-layer neural networks, take the general form of a univariate periodic function applied to an affine projection of the data. These functions have appeared in previous seminal works which demonstrate their hardness against gradient-based (Shamir'18), and Statistical Query (SQ) algorithms (Song et al.'17). We show that if (polynomially) small noise is added to the labels, the intractability of learning these functions applies to all polynomial-time algorithms, beyond gradient-based and SQ algorithms, under the aforementioned cryptographic assumptions. Moreover, we demonstrate the necessity of noise in the hardness result by designing a polynomial-time algorithm for learning certain families of such functions under exponentially small adversarial noise. Our proposed algorithm is not a gradient-based or an SQ algorithm, but is rather based on the celebrated Lenstra-Lenstra-Lov'asz (LLL) lattice basis reduction algorithm. Furthermore, in the absence of noise, this algorithm can be directly applied to solve CLWE detection (Bruna et al.'21) and phase retrieval with an optimal sample complexity of $\$d+1\$$ samples. In the former case, this improves upon the quadratic-in- $\$d\$$ sample complexity required in (Bruna et al.'21).

[PCA Initialization for Approximate Message Passing in Rotationally Invariant Models](#)

- Marco Mondelli, Ramji Venkataramanan
- abstract@[open-review](#): We study the problem of estimating a rank-1 signal in the presence of rotationally invariant noise--a class of perturbations more general than Gaussian noise. Principal Component Analysis (PCA) provides a natural estimator, and sharp results on its performance have been obtained in the high-dimensional regime. Recently, an Approximate Message Passing (AMP) algorithm has been proposed as an alternative estimator with the potential to improve the accuracy of PCA. However, the existing analysis of AMP requires an initialization that is both correlated with the signal and independent of the noise, which is often unrealistic in practice. In this work, we combine the two methods, and propose to initialize AMP with PCA. Our main result is a rigorous asymptotic characterization of the performance of this estimator. Both the AMP algorithm and its analysis differ from those previously derived in the Gaussian setting: at every iteration, our AMP algorithm requires a specific term to account for PCA initialization, while in the Gaussian case, PCA initialization affects only the first iteration of AMP. The proof is based on a two-phase artificial AMP that first approximates the PCA estimator and then mimics the true AMP. Our numerical simulations show an excellent agreement between AMP results and theoretical predictions, and suggest an interesting open direction on achieving Bayes-optimal performance.

[Automatic and Harmless Regularization with Constrained and Lexicographic Optimization: A Dynamic Barrier Approach](#)

- Chengyue Gong, Xingchao Liu, Qiang Liu
- abstract@[open-review](#): Many machine learning tasks have to make a trade-off between two loss functions, typically the main data-fitness loss and an auxiliary loss. The most widely used approach is to optimize the linear combination of the objectives, which, however, requires manual tuning of the combination coefficient and is theoretically unsuitable for non-convex functions. In this work, we consider constrained optimization as a more principled approach for trading off two losses, with a special emphasis on lexicographic optimization, a degenerated limit of constrained optimization which

optimizes a secondary loss inside the optimal set of the main loss. We propose a dynamic barrier gradient descent algorithm which provides a unified solution of both constrained and lexicographic optimization. We establish the convergence of the method for general non-convex functions.

[Corruption Robust Active Learning](#)

- Yifang Chen, Simon S. Du, Kevin G. Jamieson
- abstract@[open-review](#): We conduct theoretical studies on streaming-based active learning for binary classification under unknown adversarial label corruptions. In this setting, every time before the learner observes a sample, the adversary decides whether to corrupt the label or not. First, we show that, in a benign corruption setting (which includes the misspecification setting as a special case), with a slight enlargement on the hypothesis elimination threshold, the classical RobustCAL framework can (surprisingly) achieve nearly the same label complexity guarantee as in the non-corrupted setting. However, this algorithm can fail in the general corruption setting. To resolve this drawback, we propose a new algorithm which is provably correct without any assumptions on the presence of corruptions. Furthermore, this algorithm enjoys the minimax label complexity in the non-corrupted setting (which is achieved by RobustCAL) and only requires $\tilde{O}(C_{\text{total}})$ additional labels in the corrupted setting to achieve $\mathcal{O}(\bar{\omega} + \frac{C_{\text{total}}}{n})$, where $\bar{\omega}$ is the target accuracy, C_{total} is the total number of corruptions and n is the total number of unlabeled samples.

[Metadata-based Multi-Task Bandits with Bayesian Hierarchical Models](#)

- Runzhe Wan, Lin Ge, Rui Song
- abstract@[open-review](#): How to explore efficiently is a central problem in multi-armed bandits. In this paper, we introduce the metadata-based multi-task bandit problem, where the agent needs to solve a large number of related multi-armed bandit tasks and can leverage some task-specific features (i.e., metadata) to share knowledge across tasks. As a general framework, we propose to capture task relations through the lens of Bayesian hierarchical models, upon which a Thompson sampling algorithm is designed to efficiently learn task relations, share information, and minimize the cumulative regrets. Two concrete examples for Gaussian bandits and Bernoulli bandits are carefully analyzed. The Bayes regret for Gaussian bandits clearly demonstrates the benefits of information sharing with our algorithm. The proposed method is further supported by extensive experiments.

[Program Synthesis Guided Reinforcement Learning for Partially Observed Environments](#)

- Yichen Yang, Jeevana Priya Inala, Osbert Bastani, Yewen Pu, Armando Solar-Lezama, Martin Rinard
- abstract@[open-review](#): A key challenge for reinforcement learning is solving long-horizon planning problems. Recent work has leveraged programs to guide reinforcement learning in these settings. However, these approaches impose a high manual burden on the user since they must provide a guiding program for every new task. Partially observed environments further complicate the programming task because the program must implement a strategy that correctly, and ideally optimally, handles every possible configuration of the hidden regions of the environment. We propose a new approach, model predictive program synthesis (MPPS), that uses program synthesis to automatically generate the guiding programs. It trains a generative model to predict the unobserved portions of the world, and then synthesizes a program based on samples from this model in a way that is robust to its uncertainty. In our experiments, we show that our approach significantly outperforms non-program-guided approaches on a set of challenging benchmarks, including a 2D Minecraft-inspired environment where the agent must complete a complex sequence of subtasks to achieve its goal, and achieves a similar performance as using handcrafted programs to guide the agent. Our results demonstrate that our approach can obtain the benefits of program-guided reinforcement learning without requiring the user to provide a new guiding program for every new task.

[Robust Allocations with Diversity Constraints](#)

- Zeyu Shen, Lodewijk Gelauff, Ashish Goel, Aleksandra Korolova, Kamesh Munagala
- abstract@[open-review](#): We consider the problem of allocating divisible items among multiple agents, and consider the setting where any agent is allowed to introduce diversity constraints on the items they are allocated. We motivate this via settings where the items themselves correspond to user ad slots or task workers with attributes such as race and gender on which the principal seeks to achieve demographic parity. We consider the following question: When an agent expresses diversity constraints into an allocation rule, is the allocation of other agents hurt significantly? If this happens, the cost of introducing such constraints is disproportionately borne by agents who do not benefit from diversity. We codify this via two desiderata capturing robustness. These are no negative externality -- other agents are not hurt -- and monotonicity -- the agent enforcing the constraint does not see a large increase in value. We show in a formal sense that the Nash Welfare rule that maximizes product of agent values is uniquely positioned to be robust when diversity constraints are introduced, while almost all other natural allocation rules fail this criterion. We also show that the guarantees achieved by Nash Welfare are nearly optimal within a widely studied class of allocation rules. We finally perform an empirical simulation on real-world data that models ad allocations to show that this gap between Nash Welfare and other rules persists in the wild.

[Activation Sharing with Asymmetric Paths Solves Weight Transport Problem without Bidirectional Connection](#)

- Sunghyeon Woo, Jeongwoo Park, Jiwoo Hong, Dongsuk Jeon
- abstract@[open-review](#): One of the reasons why it is difficult for the brain to perform backpropagation (BP) is the weight transport problem, which argues forward and feedback neurons cannot share the same synaptic weights during learning in biological neural networks. Recently proposed algorithms address the weight transport problem while providing good performance similar to BP in large-scale networks. However, they require bidirectional connections between the forward and feedback neurons to train their weights, which is observed to be rare in the biological brain. In this work, we propose an Activation Sharing algorithm that removes the need for bidirectional connections between the two types of neurons. In this algorithm, hidden layer outputs (activations) are shared across multiple layers during weight updates. By applying this learning rule to both forward and feedback networks, we solve the weight transport problem without the constraint of bidirectional connections, also achieving good performance even on deep convolutional neural networks for various datasets. In addition, our algorithm could significantly reduce memory access overhead when implemented in hardware.

[BlendGAN: Implicitly GAN Blending for Arbitrary Stylized Face Generation](#)

- Mingcong Liu, Qiang Li, Zekui Qin, Guoxin Zhang, Pengfei Wan, Wen Zheng
- abstract@[open-review](#): Generative Adversarial Networks (GANs) have made a dramatic leap in high-fidelity image synthesis and stylized face generation. Recently, a layer-swapping mechanism has been developed to improve the stylization performance. However, this method is incapable of fitting arbitrary styles in a single model and requires hundreds of style-consistent training images for each style. To address the above issues, we propose BlendGAN for arbitrary stylized face generation by leveraging a flexible blending strategy and a generic artistic dataset. Specifically, we first train a self-supervised style encoder on the generic artistic dataset to extract the representations of arbitrary styles. In addition, a weighted blending module (WBM) is proposed to blend face and style representations implicitly and control the arbitrary stylization effect. By doing so, BlendGAN can gracefully fit arbitrary styles in a unified model while avoiding case-by-case preparation of style-consistent training images. To this end, we also present a novel large-scale artistic face dataset AAHQ. Extensive experiments demonstrate that BlendGAN outperforms state-of-the-art methods in terms of visual quality and style diversity for both latent-guided and reference-guided stylized face synthesis.

[Differentially Private Model Personalization](#)

- Prateek Jain, John Rush, Adam Smith, Shuang Song, Abhradeep Guha Thakurta
- abstract@[open-review](#): We study personalization of supervised learning with user-level differential privacy. Consider a setting with many users, each of whom has a training data set drawn from their own distribution $\$P_i\$$. Assuming some shared structure among the problems $\$P_i\$$, can users collectively learn the shared structure---and solve their tasks better than they could individually---while preserving the privacy of their data? We formulate this question using joint, user-level differential privacy---that is, we control what is leaked about each user's entire data set. We provide algorithms that exploit popular non-private approaches in this domain like the Almost-No-Inner-Loop (ANIL) method, and give strong user-level privacy guarantees for our general approach. When the problems $\$P_i\$$ are linear regression problems with each user's regression vector lying in a common, unknown low-dimensional subspace, we show that our efficient algorithms satisfy nearly optimal estimation error guarantees. We also establish a general, information-theoretic upper bound via an exponential mechanism-based algorithm.

[Rates of Estimation of Optimal Transport Maps using Plug-in Estimators via Barycentric Projections](#)

- NABARUN DEB, Promit Ghosal, Bodhisattva Sen
- abstract@[open-review](#): Optimal transport maps between two probability distributions $\$\mu\$$ and $\$\nu\$$ on $\$R^d\$$ have found extensive applications in both machine learning and statistics. In practice, these maps need to be estimated from data sampled according to $\$\mu\$$ and $\$\nu\$$. Plug-in estimators are perhaps most popular in estimating transport maps in the field of computational optimal transport. In this paper, we provide a comprehensive analysis of the rates of convergences for general plug-in estimators defined via barycentric projections. Our main contribution is a new stability estimate for barycentric projections which proceeds under minimal smoothness assumptions and can be used to analyze general plug-in estimators. We illustrate the usefulness of this stability estimate by first providing rates of convergence for the natural discrete-discrete and semi-discrete estimators of optimal transport maps. We then use the same stability estimate to show that, under additional smoothness assumptions of Besov type or Sobolev type, wavelet based or kernel smoothed plug-in estimators respectively speed up the rates of convergence and significantly mitigate the curse of dimensionality suffered by the natural discrete-discrete/semi-discrete estimators. As a by-product of our analysis, we also obtain faster rates of convergence for plug-in estimators of $\$W_2(\mu, \nu)\$$, the Wasserstein distance between $\$\mu\$$ and $\$\nu\$$, under the aforementioned smoothness assumptions, thereby complementing recent results in Chizat et al. (2020). Finally, we illustrate the applicability of our results in obtaining rates of convergence for Wasserstein barycenters between two probability distributions and obtaining asymptotic detection thresholds for some recent optimal-transport based tests of independence.

[Robust Generalization despite Distribution Shift via Minimum Discriminating Information](#)

- Tobias Sutter, Andreas Krause, Daniel Kuhn
- abstract@[open-review](#): Training models that perform well under distribution shifts is a central challenge in machine learning. In this paper, we introduce a modeling framework where, in addition to training data, we have partial structural knowledge of the shifted test distribution. We employ the principle of minimum discriminating information to embed the available prior knowledge, and use distributionally robust optimization to account for uncertainty due to the limited samples. By leveraging large deviation results, we obtain explicit generalization bounds with respect to the unknown shifted distribution. Lastly, we demonstrate the versatility of our framework by demonstrating it on two rather distinct applications: (1) training classifiers on systematically biased data and (2) off-policy evaluation in Markov Decision Processes.

[Soft Calibration Objectives for Neural Networks](#)

- Archit Karandikar, Nicholas Cain, Dustin Tran, Balaji Lakshminarayanan, Jonathon Shlens, Michael C. Mozer, Becca Roelofs
- abstract@[open-review](#): Optimal decision making requires that classifiers produce uncertainty estimates consistent with their empirical accuracy. However, deep neural networks are often under- or over-confident in their predictions. Consequently, methods have been developed to improve the calibration of their predictive uncertainty both during training and post-hoc. In this work, we propose differentiable losses to improve calibration based on a soft (continuous) version of the binning operation underlying popular calibration-error estimators. When incorporated into training, these soft calibration losses achieve state-of-the-art single-model ECE across multiple datasets with less than 1% decrease in accuracy. For instance, we observe an 82% reduction in ECE (70% relative to the post-hoc rescaled ECE) in exchange for a 0.7% relative decrease in accuracy relative to the cross entropy baseline on CIFAR-100. When incorporated post-training, the soft-binning-based calibration error objective improves upon temperature scaling, a popular recalibration method. Overall, experiments across losses and datasets demonstrate that using calibration-sensitive procedures yield better uncertainty estimates under dataset shift than the standard practice of using a cross entropy loss and post-hoc recalibration methods.

[Distributional Gradient Matching for Learning Uncertain Neural Dynamics Models](#)

- Lenart Treven, Philippe Wenk, Florian Dorfler, Andreas Krause
- abstract@[open-review](#): Differential equations in general and neural ODEs in particular are an essential technique in continuous-time system identification. While many deterministic learning algorithms have been designed based on numerical integration via the adjoint method, many downstream tasks such as active learning, exploration in reinforcement learning, robust control, or filtering require accurate estimates of predictive uncertainties. In this work, we propose a novel approach towards estimating epistemically uncertain neural ODEs, avoiding the numerical integration bottleneck. Instead of modeling uncertainty in the ODE parameters, we directly model uncertainties in the state space. Our algorithm distributional gradient matching (DGM) jointly trains a smoother and a dynamics model and matches their gradients via minimizing a Wasserstein loss. Our experiments show that, compared to traditional approximate inference methods based on numerical integration, our approach is faster to train, faster at predicting previously unseen trajectories, and in the context of neural ODEs, significantly more accurate.

[Shaping embodied agent behavior with activity-context priors from egocentric video](#)

- Tushar Nagarajan, Kristen Grauman
- abstract@[open-review](#): Complex physical tasks entail a sequence of object interactions, each with its own preconditions -- which can be difficult for robotic agents to learn efficiently solely through their own experience. We introduce an approach to discover activity-context priors from in-the-wild egocentric video captured with human worn cameras. For a given object, an activity-context prior represents the set of other compatible objects that are required for activities to succeed (e.g., a knife and cutting board brought together with a tomato are conducive to cutting). We encode our video-based prior as an auxiliary reward function that encourages an agent to bring compatible objects together before attempting an interaction. In this way, our model translates everyday human experience into embodied agent skills. We demonstrate our idea using egocentric EPIC-Kitchens video of people performing unscripted kitchen activities to benefit virtual household robotic agents performing various complex tasks in AI2-iTHOR, significantly accelerating agent learning.

[Adjusting for Autocorrelated Errors in Neural Networks for Time Series](#)

- Fan-Keng Sun, Chris Lang, Duane Boning
- abstract@[open-review](#): An increasing body of research focuses on using neural networks to model time series. A common assumption in training neural networks via maximum likelihood estimation on time series is that the errors across time steps are uncorrelated. However, errors are actually autocorrelated in many cases due to the temporality of the data, which makes such maximum likelihood estimations inaccurate. In this paper, in order to adjust for autocorrelated errors, we propose to learn the autocorrelation coefficient jointly with the model parameters. In our experiments, we verify the effectiveness of our approach on time series forecasting. Results across a wide range of real-world datasets with various state-of-the-art models show that our method enhances performance in almost all cases. Based on these results, we suggest empirical critical values to determine the severity of

autocorrelated errors. We also analyze several aspects of our method to demonstrate its advantages. Finally, other time series tasks are also considered to validate that our method is not restricted to only forecasting.

[A Geometric Analysis of Neural Collapse with Unconstrained Features](#)

- Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, Qing Qu
- abstract@[open-review](#): We provide the first global optimization landscape analysis of Neural Collapse -- an intriguing empirical phenomenon that arises in the last-layer classifiers and features of neural networks during the terminal phase of training. As recently reported by Papyan et al., this phenomenon implies that (i) the class means and the last-layer classifiers all collapse to the vertices of a Simplex Equiangular Tight Frame (ETF) up to scaling, and (ii) cross-example within-class variability of last-layer activations collapses to zero. We study the problem based on a simplified unconstrained feature model, which isolates the topmost layers from the classifier of the neural network. In this context, we show that the classical cross-entropy loss with weight decay has a benign global landscape, in the sense that the only global minimizers are the Simplex ETFs while all other critical points are strict saddles whose Hessian exhibit negative curvature directions. Our analysis of the simplified model not only explains what kind of features are learned in the last layer, but also shows why they can be efficiently optimized, matching the empirical observations in practical deep network architectures. These findings provide important practical implications. As an example, our experiments demonstrate that one may set the feature dimension equal to the number of classes and fix the last-layer classifier to be a Simplex ETF for network training, which reduces memory cost by over 20% on ResNet18 without sacrificing the generalization performance. The source code is available at <https://github.com/tding1/Neural-Collapse>.

[NeRS: Neural Reflectance Surfaces for Sparse-view 3D Reconstruction in the Wild](#)

- Jason Zhang, Gengshan Yang, Shubham Tulsiani, Deva Ramanan
- abstract@[open-review](#): Recent history has seen a tremendous growth of work exploring implicit representations of geometry and radiance, popularized through Neural Radiance Fields (NeRF). Such works are fundamentally based on a (implicit) {\em volumetric} representation of occupancy, allowing them to model diverse scene structure including translucent objects and atmospheric obscurants. But because the vast majority of real-world scenes are composed of well-defined surfaces, we introduce a {\em surface} analog of such implicit models called Neural Reflectance Surfaces (NeRS). NeRS learns a neural shape representation of a closed surface that is diffeomorphic to a sphere, guaranteeing water-tight reconstructions. Even more importantly, surface parameterizations allow NeRS to learn (neural) bidirectional surface reflectance functions (BRDFs) that factorize view-dependent appearance into environmental illumination, diffuse color (albedo), and specular “shininess.” Finally, rather than illustrating our results on synthetic scenes or controlled in-the-lab capture, we assemble a novel dataset of multi-view images from online marketplaces for selling goods. Such “in-the-wild” multi-view image sets pose a number of challenges, including a small number of views with unknown/rough camera estimates. We demonstrate that surface-based neural reconstructions enable learning from such data, outperforming volumetric neural rendering-based reconstructions. We hope that NeRS serves as a first step toward building scalable, high-quality libraries of real-world shape, materials, and illumination.

[Unleashing the Power of Contrastive Self-Supervised Visual Models via Contrast-Regularized Fine-Tuning](#)

- Yifan Zhang, Bryan Hooi, Dapeng Hu, Jian Liang, Jiashi Feng
- abstract@[open-review](#): Contrastive self-supervised learning (CSL) has attracted increasing attention for model pre-training via unlabeled data. The resulted CSL models provide instance-discriminative visual features that are uniformly scattered in the feature space. During deployment, the common practice is to directly fine-tune CSL models with cross-entropy, which however may not be the best strategy in practice. Although cross-entropy tends to separate inter-class features, the resulting models still have limited capability for reducing intra-class feature scattering that exists in CSL models. In this paper, we investigate whether applying contrastive learning to fine-tuning would bring further benefits, and analytically find that optimizing the contrastive loss benefits both discriminative representation learning and model optimization during fine-tuning. Inspired by these findings, we propose Contrast-regularized tuning (Core-tuning), a new approach for fine-tuning CSL models. Instead of simply adding the contrastive loss to the objective of fine-tuning, Core-tuning further applies a novel hard pair mining strategy for more effective contrastive fine-tuning, as well as smoothing the decision boundary to better exploit the learned discriminative feature space. Extensive experiments on image classification and semantic segmentation verify the effectiveness of Core-tuning.

[Discovery of Options via Meta-Learned Subgoals](#)

- Vivek Veeriah, Tom Zahavy, Matteo Hessel, Zhongwen Xu, Junhyuk Oh, Iurii Kemaev, Hado P. van Hasselt, David Silver, Satinder Singh
- abstract@[open-review](#): Temporal abstractions in the form of options have been shown to help reinforcement learning (RL) agents learn faster. However, despite prior work on this topic, the problem of discovering options through interaction with an environment remains a challenge. In this paper, we introduce a novel meta-gradient approach for discovering useful options in multi-task RL environments. Our approach is based on a manager-worker decomposition of the RL agent, in which a manager maximises rewards from the environment by learning a task-dependent policy over both a set of task-independent discovered-options and primitive actions. The option-reward and termination functions that define a subgoal for each option are parameterised as neural networks and trained via meta-gradients to maximise their usefulness. Empirical analysis on gridworld and DeepMind Lab tasks show that: (1) our approach can discover meaningful and diverse temporally-extended options in multi-task RL domains, (2) the discovered options are frequently used by the agent while learning to solve the training tasks, and (3) that the discovered options help a randomly initialised manager learn faster in completely new tasks.

[Near-Optimal Lower Bounds For Convex Optimization For All Orders of Smoothness](#)

- Ankit Garg, Robin Kothari, Praneeth Netrapalli, Suhail Sherif
- abstract@[open-review](#): We study the complexity of optimizing highly smooth convex functions. For a positive integer \$p\$, we want to find an \$\epsilon\$-approximate minimum of a convex function \$f\$, given oracle access to the function and its first \$p\$ derivatives, assuming that the \$p\$th derivative of \$f\$ is Lipschitz. Recently, three independent research groups (Jiang et al., PLMR 2019; Gasnikov et al., PLMR 2019; Bubeck et al., PLMR 2019) developed a new algorithm that solves this problem with \$\tilde{O}\left(1/\epsilon^{\frac{2}{3p+1}}\right)\$ oracle calls for constant \$p\$. This is known to be optimal (up to log factors) for deterministic algorithms, but known lower bounds for randomized algorithms do not match this bound. We prove a new lower bound that matches this bound (up to log factors), and holds not only for randomized algorithms, but also for quantum algorithms.

[Topology-Imbalance Learning for Semi-Supervised Node Classification](#)

- Deli Chen, Yankai Lin, Guangxiang Zhao, Xuancheng Ren, Peng Li, Jie Zhou, Xu Sun
- abstract@[open-review](#): The class imbalance problem, as an important issue in learning node representations, has drawn increasing attention from the community. Although the imbalance considered by existing studies roots from the unequal quantity of labeled examples in different classes (quantity imbalance), we argue that graph data expose a unique source of imbalance from the asymmetric topological properties of the labeled nodes, i.e., labeled nodes are not equal in terms of their structural role in the graph (topology imbalance). In this work, we first probe the previously unknown topology-imbalance issue, including its characteristics, causes, and threats to semisupervised node classification learning. We then provide a unified view to jointly analyzing the quantity- and topology- imbalance issues by considering the node influence shift phenomenon with the Label Propagation algorithm. In light of our analysis, we devise an influence conflict detection-based metric Totoro to measure the degree of graph topology imbalance and propose a model-agnostic method ReNode to address the topology-imbalance issue by re-weighting the influence of labeled nodes adaptively based on their relative positions to class boundaries. Systematic experiments demonstrate the effectiveness and generalizability of our method in relieving topology-imbalance

issue and promoting semi-supervised node classification. The further analysis unveils varied sensitivity of different graph neural networks (GNNs) to topology imbalance, which may serve as a new perspective in evaluating GNN architectures.

[Gradient Inversion with Generative Image Prior](#)

- Jinwoo Jeon, jaechang Kim, Kangwook Lee, Sewoong Oh, Jungseul Ok
- abstract@[open-review](#): Federated Learning (FL) is a distributed learning framework, in which the local data never leaves clients' devices to preserve privacy, and the server trains models on the data via accessing only the gradients of those local data. Without further privacy mechanisms such as differential privacy, this leaves the system vulnerable against an attacker who inverts those gradients to reveal clients' sensitive data. However, a gradient is often insufficient to reconstruct the user data without any prior knowledge. By exploiting a generative model pretrained on the data distribution, we demonstrate that data privacy can be easily breached. Further, when such prior knowledge is unavailable, we investigate the possibility of learning the prior from a sequence of gradients seen in the process of FL training. We experimentally show that the prior in a form of generative model is learnable from iterative interactions in FL. Our findings demonstrate that additional mechanisms are necessary to prevent privacy leakage in FL.

[Beta-CROWN: Efficient Bound Propagation with Per-neuron Split Constraints for Neural Network Robustness Verification](#)

- Shiqi Wang, Huan Zhang, Kaidi Xu, Xue Lin, Suman Jana, Cho-Jui Hsieh, J. Zico Kolter
- abstract@[open-review](#): Bound propagation based incomplete neural network verifiers such as CROWN are very efficient and can significantly accelerate branch-and-bound (BaB) based complete verification of neural networks. However, bound propagation cannot fully handle the neuron split constraints introduced by BaB commonly handled by expensive linear programming (LP) solvers, leading to loose bounds and hurting verification efficiency. In this work, we develop β -CROWN, a new bound propagation based method that can fully encode neuron splits via optimizable parameters β constructed from either primal or dual space. When jointly optimized in intermediate layers, β -CROWN generally produces better bounds than typical LP verifiers with neuron split constraints, while being as efficient and parallelizable as CROWN on GPUs. Applied to complete robustness verification benchmarks, β -CROWN with BaB is up to three orders of magnitude faster than LP-based BaB methods, and is notably faster than all existing approaches while producing lower timeout rates. By terminating BaB early, our method can also be used for efficient incomplete verification. We consistently achieve higher verified accuracy in many settings compared to powerful incomplete verifiers, including those based on convex barrier breaking techniques. Compared to the typically tightest but very costly semidefinite programming (SDP) based incomplete verifiers, we obtain higher verified accuracy with three orders of magnitudes less verification time. Our algorithm empowered the α, β -CROWN (α - β -CROWN) verifier, the winning tool in VNN-COMP 2021. Our code is available at <http://PaperCode.cc/BetaCROWN>.

[Autobahn: Automorphism-based Graph Neural Nets](#)

- Erik Thiede, Wenda Zhou, Risi Kondor
- abstract@[open-review](#): We introduce Automorphism-based graph neural networks (Autobahn), a new family of graph neural networks. In an Autobahn, we decompose the graph into a collection of subgraphs and apply local convolutions that are equivariant to each subgraph's automorphism group. Specific choices of local neighborhoods and subgraphs recover existing architectures such as message passing neural networks. Our formalism also encompasses novel architectures: as an example, we introduce a graph neural network that decomposes the graph into paths and cycles. The resulting convolutions reflect the natural way that parts of the graph can transform, preserving the intuitive meaning of convolution without sacrificing global permutation equivariance. We validate our approach by applying Autobahn to molecular graphs, where it achieves results competitive with state-of-the-art message passing algorithms.

[Data Augmentation Can Improve Robustness](#)

- Sylvestre-Alvise Rebuffi, Sven Gowal, Dan Andrei Calian, Florian Stimberg, Olivia Wiles, Timothy A Mann
- abstract@[open-review](#): Adversarial training suffers from robust overfitting, a phenomenon where the robust test accuracy starts to decrease during training. In this paper, we focus on reducing robust overfitting by using common data augmentation schemes. We demonstrate that, contrary to previous findings, when combined with model weight averaging, data augmentation can significantly boost robust accuracy. Furthermore, we compare various augmentations techniques and observe that spatial composition techniques work the best for adversarial training. Finally, we evaluate our approach on CIFAR-10 against ℓ_∞ and ℓ_2 norm-bounded perturbations of size $\epsilon = 8/255$ and $\epsilon = 128/255$, respectively. We show large absolute improvements of +2.93% and +2.16% in robust accuracy compared to previous state-of-the-art methods. In particular, against ℓ_∞ norm-bounded perturbations of size $\epsilon = 8/255$, our model reaches 60.07% robust accuracy without using any external data. We also achieve a significant performance boost with this approach while using other architectures and datasets such as CIFAR-100, SVHN and TinyImageNet.

[Deep Explicit Duration Switching Models for Time Series](#)

- Abdul Fatir Ansari, Konstantinos Benidis, Richard Kurle, Ali Caner Turkmen, Harold Soh, Alexander J. Smola, Bernie Wang, Tim Januschowski
- abstract@[open-review](#): Many complex time series can be effectively subdivided into distinct regimes that exhibit persistent dynamics. Discovering the switching behavior and the statistical patterns in these regimes is important for understanding the underlying dynamical system. We propose the Recurrent Explicit Duration Switching Dynamical System (RED-SDS), a flexible model that is capable of identifying both state- and time-dependent switching dynamics. State-dependent switching is enabled by a recurrent state-to-switch connection and an explicit duration count variable is used to improve the time-dependent switching behavior. We demonstrate how to perform efficient inference using a hybrid algorithm that approximates the posterior of the continuous states via an inference network and performs exact inference for the discrete switches and counts. The model is trained by maximizing a Monte Carlo lower bound of the marginal log-likelihood that can be computed efficiently as a byproduct of the inference routine. Empirical results on multiple datasets demonstrate that RED-SDS achieves considerable improvement in time series segmentation and competitive forecasting performance against the state of the art.

[Shared Independent Component Analysis for Multi-Subject Neuroimaging](#)

- Hugo Richard, Pierre Ablin, Bertrand Thirion, Alexandre Gramfort, Aapo Hyvärinen
- abstract@[open-review](#): We consider shared response modeling, a multi-view learning problem where one wants to identify common components from multiple datasets or views. We introduce Shared Independent Component Analysis (ShICA) that models eachview as a linear transform of shared independent components contaminated by additive Gaussian noise. We show that this model is identifiable if the components are either non-Gaussian or have enough diversity in noise variances. We then show that in some cases multi-set canonical correlation analysis can recover the correct unmixing matrices, but that even a small amount of sampling noise makes Multiset CCA fail. To solve this problem, we propose to use joint diagonalization after Multiset CCA, leading to a new approach called ShICA-J. We show via simulations that ShICA-J leads to improved results while being very fast to fit. While ShICA-J is based on second-order statistics, we further propose to leverage non-Gaussianity of the components using a maximum-likelihood method, ShICA-ML, that is both more accurate and more costly. Further, ShICA comes with a principled method for shared components estimation. Finally, we provide empirical evidence on fMRI and MEG datasets that ShICA yields more accurate estimation of the components than alternatives.

[Shape from Blur: Recovering Textured 3D Shape and Motion of Fast Moving Objects](#)

- Denys Rozumnyi, Martin R. Oswald, Vittorio Ferrari, Marc Pollefeys
- abstract@[open-review](#): We address the novel task of jointly reconstructing the 3D shape, texture, and motion of an object from a single motion-blurred image. While previous approaches address the deblurring problem only in the 2D image domain, our proposed rigorous modeling of all object properties in the 3D domain enables the correct description of arbitrary object motion. This leads to significantly better image decomposition and sharper deblurring results. We model the observed appearance of a motion-blurred object as a combination of the background and a 3D object with constant translation and rotation. Our method minimizes a loss on reconstructing the input image via differentiable rendering with suitable regularizers. This enables estimating the textured 3D mesh of the blurred object with high fidelity. Our method substantially outperforms competing approaches on several benchmarks for fast moving objects deblurring. Qualitative results show that the reconstructed 3D mesh generates high-quality temporal super-resolution and novel views of the deblurred object.

[Batched Thompson Sampling](#)

- Cem Kalkanli, Ayfer Ozgur
- abstract@[open-review](#): We introduce a novel anytime batched Thompson sampling policy for multi-armed bandits where the agent observes the rewards of her actions and adjusts her policy only at the end of a small number of batches. We show that this policy simultaneously achieves a problem dependent regret of order $\mathcal{O}(\log(T))$ and a minimax regret of order $\mathcal{O}(\sqrt{T \log(T)})$ while the number of batches can be bounded by $\mathcal{O}(\log(T))$ independent of the problem instance over a time horizon T . We also prove that in expectation the instance dependent batch complexity of our policy is of order $\mathcal{O}(\log \log(T))$. These results indicate that Thompson sampling performs competitively with recently proposed algorithms for the batched setting, which optimize the batch structure for a given time horizon T and prioritize exploration in the beginning of the experiment to eliminate suboptimal actions. Unlike these algorithms, the batched Thompson sampling algorithm we propose is an anytime policy, i.e. it operates without the knowledge of the time horizon T , and as such it is the only anytime algorithm that achieves optimal regret with $\mathcal{O}(\log \log(T))$ expected batch complexity. This is achieved through a dynamic batching strategy, which uses the agents estimates to adaptively increase the batch duration.

[Delayed Gradient Averaging: Tolerate the Communication Latency for Federated Learning](#)

- Ligeng Zhu, Hongzhou Lin, Yao Lu, Yujun Lin, Song Han
- abstract@[open-review](#): Federated Learning is an emerging direction in distributed machine learning that enables jointly training a model without sharing the data. Since the data is distributed across many edge devices through wireless / long-distance connections, federated learning suffers from inevitable high communication latency. However, the latency issues are undermined in the current literature [15] and existing approaches such as FedAvg [27] become less efficient when the latency increases. To overcome the problem, we propose \textbf{Delayed Gradient Averaging} (DGA), which delays the averaging step to improve efficiency and allows local computation in parallel to communication. We theoretically prove that DGA attains a similar convergence rate as FedAvg, and empirically show that our algorithm can tolerate high network latency without compromising accuracy. Specifically, we benchmark the training speed on various vision (CIFAR, ImageNet) and language tasks (Shakespeare), with both IID and non-IID partitions, and show DGA can bring 2.55\$ times \$ to 4.07\$ times \$ speedup. Moreover, we built a 16-node Raspberry Pi cluster and show that DGA can consistently speed up real-world federated learning applications.

[Focal Attention for Long-Range Interactions in Vision Transformers](#)

- Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Xiyang Dai, Bin Xiao, Lu Yuan, Jianfeng Gao
- abstract@[open-review](#): Recently, Vision Transformer and its variants have shown great promise on various computer vision tasks. The ability to capture local and global visual dependencies through self-attention is the key to its success. But it also brings challenges due to quadratic computational overhead, especially for the high-resolution vision tasks (e.g., object detection). Many recent works have attempted to reduce the cost and improve model performance by applying either coarse-grained global attention or fine-grained local attention. However, both approaches cripple the modeling power of the original self-attention mechanism of multi-layer Transformers, leading to sub-optimal solutions. In this paper, we present focal attention, a new attention mechanism that incorporates both fine-grained local and coarse-grained global interactions. In this new mechanism, each token attends its closest surrounding tokens at the fine granularity and the tokens far away at a coarse granularity and thus can capture both short- and long-range visual dependencies efficiently and effectively. With focal attention, we propose a new variant of Vision Transformer models, called Focal Transformers, which achieve superior performance over the state-of-the-art (SoTA) Vision Transformers on a range of public image classification and object detection benchmarks. In particular, our Focal Transformer models with a moderate size of 51.1M and a large size of 89.8M achieve 83.6% and 84.0% Top-1 accuracy, respectively, on ImageNet classification at 224×224. When employed as the backbones, Focal Transformers achieve consistent and substantial improvements over the current SoTA Swin Transformers [44] across 6 different object detection methods. Our largest Focal Transformer yields 58.7/59.0 box mAPs and 50.9/51.3 mask mAPs on COCO mini-val/test-dev, and 55.4 mIoU on ADE20K for semantic segmentation, creating new SoTA on three of the most challenging computer vision tasks.

[Scalable and Stable Surrogates for Flexible Classifiers with Fairness Constraints](#)

- Henry C Bendekgey, Erik Sudderth
- abstract@[open-review](#): We investigate how fairness relaxations scale to flexible classifiers like deep neural networks for images and text. We analyze an easy-to-use and robust way of imposing fairness constraints when training, and through this framework prove that some prior fairness surrogates exhibit degeneracies for non-convex models. We resolve these problems via three new surrogates: an adaptive data re-weighting, and two smooth upper-bounds that are provably more robust than some previous methods. Our surrogates perform comparably to the state-of-the-art on low-dimensional fairness benchmarks, while achieving superior accuracy and stability for more complex computer vision and natural language processing tasks.

[Residual Pathway Priors for Soft Equivariance Constraints](#)

- Marc Finzi, Gregory Benton, Andrew G. Wilson
- abstract@[open-review](#): Models such as convolutional neural networks restrict the hypothesis space to a set of functions satisfying equivariance constraints, and improve generalization in problems by capturing relevant symmetries. However, symmetries are often only partially respected, preventing models with restriction biases from fitting the data. We introduce Residual Pathway Priors (RPPs) as a method for converting hard architectural constraints into soft priors, guiding models towards structured solutions while retaining the ability to capture additional complexity. RPPs are resilient to approximate or misspecified symmetries, and are as effective as fully constrained models even when symmetries are exact. We show that RPPs provide compelling performance on both model-free and model-based reinforcement learning problems, where contact forces and directional rewards violate the assumptions of equivariant networks. Finally, we demonstrate that RPPs have broad applicability, including dynamical systems, regression, and classification.

[Optimal Algorithms for Stochastic Contextual Preference Bandits](#)

- Aadirupa Saha
- abstract@[open-review](#): We consider the problem of preference bandits in the contextual setting. At each round, the learner is presented with a context set of K items, chosen randomly from a potentially infinite set of arms $\mathcal{D} \subseteq \mathbf{R}^d$. However, unlike classical contextual bandits, our framework only allows the learner to receive feedback in terms of item preferences: At each round, the learner is allowed to play a subset of size q (any $q \in \{2, \dots, K\}$) upon which only a (noisy) winner of the subset is revealed. Yet, same as the classical setup, the goal is still to compete against the best context arm at each round. The problem is relevant in various online decision-making scenarios, including recommender systems, information

retrieval, tournament ranking--typically any application where it's easier to elicit the items' relative strength instead of their absolute scores. To the best of our knowledge, this work is the first to consider preference-based stochastic contextual bandits for potentially infinite decision spaces. We start with presenting two algorithms for the special case of pairwise preferences ($q=2$): The first algorithm is simple and easy to implement with an $\tilde{O}(d\sqrt{T})$ regret guarantee, while the second algorithm is shown to achieve the optimal $\tilde{O}(\sqrt{dT})$ regret, as follows from our $\Omega(\sqrt{dT})$ matching lower bound analysis. We then proceed to analyze the problem for any general q -subsetwise preferences ($q \geq 2$), where surprisingly, our lower bound proves the fundamental performance limit to be $\Omega(\sqrt{dT})$ yet again, independent of the subset size q . Following this, we propose a matching upper bound algorithm justifying the tightness of our results. This implies having access to subsetwise preferences does not help in faster information aggregation for our feedback model. All the results are corroborated empirically against existing baselines.

Tight High Probability Bounds for Linear Stochastic Approximation with Fixed Stepsize

- Alain Durmus, Eric Moulines, Alexey Naumov, Sergey Samsonov, Kevin Scaman, Hoi-To Wai
- abstract@[open-review](#): This paper provides a non-asymptotic analysis of linear stochastic approximation (LSA) algorithms with fixed stepsize. This family of methods arises in many machine learning tasks and is used to obtain approximate solutions of a linear system $\bar{A}\theta = \bar{b}$ for which \bar{A} and \bar{b} can only be accessed through random estimates $\{\langle \bar{A}_n, \bar{b}_n \rangle : n \in \mathbb{N}^*\}$. Our analysis is based on new results regarding moments and high probability bounds for products of matrices which are shown to be tight. We derive high probability bounds on the performance of LSA under weaker conditions on the sequence $\{\langle \bar{A}_n, \bar{b}_n \rangle : n \in \mathbb{N}^*\}$ than previous works. However, in contrast, we establish polynomial concentration bounds with order depending on the stepsize. We show that our conclusions cannot be improved without additional assumptions on the sequence of random matrices $\{\bar{A}_n : n \in \mathbb{N}^*\}$, and in particular that no Gaussian or exponential high probability bounds can hold. Finally, we pay a particular attention to establishing bounds with sharp order with respect to the number of iterations and the stepsize and whose leading terms contain the covariance matrices appearing in the central limit theorems.

Learning Large Neighborhood Search Policy for Integer Programming

- Yaxin Wu, Wen Song, Zhiguang Cao, Jie Zhang
- abstract@[open-review](#): We propose a deep reinforcement learning (RL) method to learn large neighborhood search (LNS) policy for integer programming (IP). The RL policy is trained as the destroy operator to select a subset of variables at each step, which is reoptimized by an IP solver as the repair operator. However, the combinatorial number of variable subsets prevents direct application of typical RL algorithms. To tackle this challenge, we represent all subsets by factorizing them into binary decisions on each variable. We then design a neural network to learn policies for each variable in parallel, trained by a customized actor-critic algorithm. We evaluate the proposed method on four representative IP problems. Results show that it can find better solutions than SCIP in much less time, and significantly outperform other LNS baselines with the same runtime. Moreover, these advantages notably persist when the policies generalize to larger problems. Further experiments with Gurobi also reveal that our method can outperform this state-of-the-art commercial solver within the same time limit.

Dynamic Trace Estimation

- Prathamesh Dharangutte, Christopher Musco
- abstract@[open-review](#): We study a dynamic version of the implicit trace estimation problem. Given access to an oracle for computing matrix-vector multiplications with a dynamically changing matrix A , our goal is to maintain an accurate approximation to A 's trace using as few multiplications as possible. We present a practical algorithm for solving this problem and prove that, in a natural setting, its complexity is quadratically better than the standard solution of repeatedly applying Hutchinson's stochastic trace estimator. We also provide an improved algorithm assuming additional common assumptions on A 's dynamic updates. We support our theory with empirical results, showing significant computational improvements on three applications in machine learning and network science: tracking moments of the Hessian spectral density during neural network optimization, counting triangles and estimating natural connectivity in a dynamically changing graph.

Provable Representation Learning for Imitation with Contrastive Fourier Features

- Ofir Nachum, Mengjiao Yang
- abstract@[open-review](#): In imitation learning, it is common to learn a behavior policy to match an unknown target policy via max-likelihood training on a collected set of target demonstrations. In this work, we consider using offline experience datasets -- potentially far from the target distribution -- to learn low-dimensional state representations that provably accelerate the sample-efficiency of downstream imitation learning. A central challenge in this setting is that the unknown target policy itself may not exhibit low-dimensional behavior, and so there is a potential for the representation learning objective to alias states in which the target policy acts differently. Circumventing this challenge, we derive a representation learning objective that provides an upper bound on the performance difference between the target policy and a low-dimensional policy trained with max-likelihood, and this bound is tight regardless of whether the target policy itself exhibits low-dimensional structure. Moving to the practicality of our method, we show that our objective can be implemented as contrastive learning, in which the transition dynamics are approximated by either an implicit energy-based model or, in some special cases, an implicit linear model with representations given by random Fourier features. Experiments on both tabular environments and high-dimensional Atari games provide quantitative evidence for the practical benefits of our proposed objective.

MICo: Improved representations via sampling-based state similarity for Markov decision processes

- Pablo Samuel Castro, Tyler Kastner, Prakash Panangaden, Mark Rowland
- abstract@[open-review](#): We present a new behavioural distance over the state space of a Markov decision process, and demonstrate the use of this distance as an effective means of shaping the learnt representations of deep reinforcement learning agents. While existing notions of state similarity are typically difficult to learn at scale due to high computational cost and lack of sample-based algorithms, our newly-proposed distance addresses both of these issues. In addition to providing detailed theoretical analyses, we provide empirical evidence that learning this distance alongside the value function yields structured and informative representations, including strong results on the Arcade Learning Environment benchmark.

Counterfactual Explanations in Sequential Decision Making Under Uncertainty

- Stratis Tsirtsis, Abir De, Manuel Rodriguez
- abstract@[open-review](#): Methods to find counterfactual explanations have predominantly focused on one-step decision making processes. In this work, we initiate the development of methods to find counterfactual explanations for decision making processes in which multiple, dependent actions are taken sequentially over time. We start by formally characterizing a sequence of actions and states using finite horizon Markov decision processes and the Gumbel-Max structural causal model. Building upon this characterization, we formally state the problem of finding counterfactual explanations for sequential decision making processes. In our problem formulation, the counterfactual explanation specifies an alternative sequence of actions differing in at most k actions from the observed sequence that could have led the observed process realization to a better outcome. Then, we introduce a polynomial time algorithm based on dynamic programming to build a counterfactual policy that is guaranteed to always provide the optimal counterfactual explanation on every possible realization of the counterfactual environment dynamics. We validate our algorithm using both synthetic and real data from cognitive behavioral therapy and show that the counterfactual explanations our algorithm finds can provide valuable insights to enhance sequential decision making under uncertainty.

[Streaming Linear System Identification with Reverse Experience Replay](#)

- Suhas Kowshik, Dheeraj Nagaraj, Prateek Jain, Praneeth Netrapalli
- abstract@[open-review](#): We consider the problem of estimating a linear time-invariant (LTI) dynamical system from a single trajectory via streaming algorithms, which is encountered in several applications including reinforcement learning (RL) and time-series analysis. While the LTI system estimation problem is well-studied in the {em offline} setting, the practically important streaming/online setting has received little attention. Standard streaming methods like stochastic gradient descent (SGD) are unlikely to work since streaming points can be highly correlated. In this work, we propose a novel streaming algorithm, SGD with Reverse Experience Replay (SGD-RER), that is inspired by the experience replay (ER) technique popular in the RL literature. SGD-RER divides data into small buffers and runs SGD backwards on the data stored in the individual buffers. We show that this algorithm exactly deconstructs the dependency structure and obtains information theoretically optimal guarantees for both parameter error and prediction error. Thus, we provide the first -- to the best of our knowledge -- optimal SGD-style algorithm for the classical problem of linear system identification with a first order oracle. Furthermore, SGD-RER can be applied to more general settings like sparse LTI identification with known sparsity pattern, and non-linear dynamical systems. Our work demonstrates that the knowledge of data dependency structure can aid us in designing statistically and computationally efficient algorithms which can ``decorrelate'' streaming samples.

[SmoothMix: Training Confidence-calibrated Smoothed Classifiers for Certified Robustness](#)

- Jongheon Jeong, Sejun Park, Minkyu Kim, Heung-Chang Lee, Do-Guk Kim, Jinwoo Shin
- abstract@[open-review](#): Randomized smoothing is currently a state-of-the-art method to construct a certifiably robust classifier from neural networks against ℓ_2 -adversarial perturbations. Under the paradigm, the robustness of a classifier is aligned with the prediction confidence, i.e., the higher confidence from a smoothed classifier implies the better robustness. This motivates us to rethink the fundamental trade-off between accuracy and robustness in terms of calibrating confidences of a smoothed classifier. In this paper, we propose a simple training scheme, coined SmoothMix, to control the robustness of smoothed classifiers via self-mixup: it trains on convex combinations of samples along the direction of adversarial perturbation for each input. The proposed procedure effectively identifies over-confident, near off-class samples as a cause of limited robustness in case of smoothed classifiers, and offers an intuitive way to adaptively set a new decision boundary between these samples for better robustness. Our experimental results demonstrate that the proposed method can significantly improve the certified ℓ_2 -robustness of smoothed classifiers compared to existing state-of-the-art robust training methods.

[Action-guided 3D Human Motion Prediction](#)

- Jiangxin Sun, Zihang Lin, Xintong Han, Jian-Fang Hu, Jia Xu, Wei-Shi Zheng
- abstract@[open-review](#): The ability of forecasting future human motion is important for human-machine interaction systems to understand human behaviors and make interaction. In this work, we focus on developing models to predict future human motion from past observed video frames. Motivated by the observation that human motion is closely related to the action being performed, we propose to explore action context to guide motion prediction. Specifically, we construct an action-specific memory bank to store representative motion dynamics for each action category, and design a query-read process to retrieve some motion dynamics from the memory bank. The retrieved dynamics are consistent with the action depicted in the observed video frames and serve as a strong prior knowledge to guide motion prediction. We further formulate an action constraint loss to ensure the global semantic consistency of the predicted motion. Extensive experiments demonstrate the effectiveness of the proposed approach, and we achieve state-of-the-art performance on 3D human motion prediction.

[Meta-Learning the Search Distribution of Black-Box Random Search Based Adversarial Attacks](#)

- Maksym Yatsura, Jan Metzen, Matthias Hein
- abstract@[open-review](#): Adversarial attacks based on randomized search schemes have obtained state-of-the-art results in black-box robustness evaluation recently. However, as we demonstrate in this work, their efficiency in different query budget regimes depends on manual design and heuristic tuning of the underlying proposal distributions. We study how this issue can be addressed by adapting the proposal distribution online based on the information obtained during the attack. We consider Square Attack, which is a state-of-the-art score-based black-box attack, and demonstrate how its performance can be improved by a learned controller that adjusts the parameters of the proposal distribution online during the attack. We train the controller using gradient-based end-to-end training on a CIFAR10 model with white box access. We demonstrate that plugging the learned controller into the attack consistently improves its black-box robustness estimate in different query regimes by up to 20% for a wide range of different models with black-box access. We further show that the learned adaptation principle transfers well to the other data distributions such as CIFAR100 or ImageNet and to the targeted attack setting.

[Validating the Lottery Ticket Hypothesis with Inertial Manifold Theory](#)

- Zeru Zhang, Jiayin Jin, Zijie Zhang, Yang Zhou, Xin Zhao, Jiaxiang Ren, Ji Liu, Lingfei Wu, Ruoming Jin, Dejing Dou
- abstract@[open-review](#): Despite achieving remarkable efficiency, traditional network pruning techniques often follow manually-crafted heuristics to generate pruned sparse networks. Such heuristic pruning strategies are hard to guarantee that the pruned networks achieve test accuracy comparable to the original dense ones. Recent works have empirically identified and verified the Lottery Ticket Hypothesis (LTH): a randomly-initialized dense neural network contains an extremely sparse subnetwork, which can be trained to achieve similar accuracy to the former. Due to the lack of theoretical evidence, they often need to run multiple rounds of expensive training and pruning over the original large networks to discover the sparse subnetworks with low accuracy loss. By leveraging dynamical systems theory and inertial manifold theory, this work theoretically verifies the validity of the LTH. We explore the possibility of theoretically lossless pruning as well as one-time pruning, compared with existing neural network pruning and LTH techniques. We reformulate the neural network optimization problem as a gradient dynamical system and reduce this high-dimensional system onto inertial manifolds to obtain a low-dimensional system regarding pruned subnetworks. We demonstrate the precondition and existence of pruned subnetworks and prune the original networks in terms of the gap in their spectrum that make the subnetworks have the smallest dimensions.

[Are My Deep Learning Systems Fair? An Empirical Study of Fixed-Seed Training](#)

- Shangshu Qian, Viet Hung Pham, Thibaud Lutellier, Zeou Hu, Jungwon Kim, Lin Tan, Yaoliang Yu, Jiahao Chen, Sameena Shah
- abstract@[open-review](#): Deep learning (DL) systems have been gaining popularity in critical tasks such as credit evaluation and crime prediction. Such systems demand fairness. Recent work shows that DL software implementations introduce variance: identical DL training runs (i.e., identical network, data, configuration, software, and hardware) with a fixed seed produce different models. Such variance could make DL models and networks violate fairness compliance laws, resulting in negative social impact. In this paper, we conduct the first empirical study to quantify the impact of software implementation on the fairness and its variance of DL systems. Our study of 22 mitigation techniques and five baselines reveals up to 12.6% fairness variance across identical training runs with identical seeds. In addition, most debiasing algorithms have a negative impact on the model such as reducing model accuracy, increasing fairness variance, or increasing accuracy variance. Our literature survey shows that while fairness is gaining popularity in artificial intelligence (AI) related conferences, only 34.4% of the papers use multiple identical training runs to evaluate their approach, raising concerns about their results' validity. We call for better fairness evaluation and testing protocols to improve fairness and fairness variance of DL systems as well as DL research validity and reproducibility at large.

[Rectangular Flows for Manifold Learning](#)

- Anthony L. Caterini, Gabriel Loaiza-Ganem, Geoff Pleiss, John P. Cunningham
- abstract@[open-review](#): Normalizing flows are invertible neural networks with tractable change-of-volume terms, which allow optimization of their parameters to be efficiently performed via maximum likelihood. However, data of interest are typically assumed to live in some (often unknown) low-dimensional manifold embedded in a high-dimensional ambient space. The result is a modelling mismatch since -- by construction -- the invertibility requirement implies high-dimensional support of the learned distribution. Injective flows, mappings from low- to high-dimensional spaces, aim to fix this discrepancy by learning distributions on manifolds, but the resulting volume-change term becomes more challenging to evaluate. Current approaches either avoid computing this term entirely using various heuristics, or assume the manifold is known beforehand and therefore are not widely applicable. Instead, we propose two methods to tractably calculate the gradient of this term with respect to the parameters of the model, relying on careful use of automatic differentiation and techniques from numerical linear algebra. Both approaches perform end-to-end nonlinear manifold learning and density estimation for data projected onto this manifold. We study the trade-offs between our proposed methods, empirically verify that we outperform approaches ignoring the volume-change term by more accurately learning manifolds and the corresponding distributions on them, and show promising results on out-of-distribution detection. Our code is available at <https://github.com/layer6ai-labs/rectangular-flows>.

[On the Generative Utility of Cyclic Conditionals](#)

- Chang Liu, Haoyue Tang, Tao Qin, Jintao Wang, Tie-Yan Liu
- abstract@[open-review](#): We study whether and how can we model a joint distribution $p(x,z)$ using two conditional models $p(x|z)$ and $q(z|x)$ that form a cycle. This is motivated by the observation that deep generative models, in addition to a likelihood model $p(x|z)$, often also use an inference model $q(z|x)$ for extracting representation, but they rely on a usually uninformative prior distribution $p(z)$ to define a joint distribution, which may render problems like posterior collapse and manifold mismatch. To explore the possibility to model a joint distribution using only $p(x|z)$ and $q(z|x)$, we study their compatibility and determinacy, corresponding to the existence and uniqueness of a joint distribution whose conditional distributions coincide with them. We develop a general theory for operable equivalence criteria for compatibility, and sufficient conditions for determinacy. Based on the theory, we propose a novel generative modeling framework CyGen that only uses the two cyclic conditional models. We develop methods to achieve compatibility and determinacy, and to use the conditional models to fit and generate data. With the prior constraint removed, CyGen better fits data and captures more representative features, supported by both synthetic and real-world experiments.

[Structural Credit Assignment in Neural Networks using Reinforcement Learning](#)

- Dhawal Gupta, Gabor Mihucz, Matthew Schlegel, James Kostas, Philip S. Thomas, Martha White
- abstract@[open-review](#): Structural credit assignment in neural networks is a long-standing problem, with a variety of alternatives to backpropagation proposed to allow for local training of nodes. One of the early strategies was to treat each node as an agent and use a reinforcement learning method called REINFORCE to update each node locally with only a global reward signal. In this work, we revisit this approach and investigate if we can leverage other reinforcement learning approaches to improve learning. We first formalize training a neural network as a finite-horizon reinforcement learning problem and discuss how this facilitates using ideas from reinforcement learning like off-policy learning. We show that the standard on-policy REINFORCE approach, even with a variety of variance reduction approaches, learns suboptimal solutions. We introduce an off-policy approach, to facilitate reasoning about the greedy action for other agents and help overcome stochasticity in other agents. We conclude by showing that these networks of agents can be more robust to correlated samples when learning online.

[A Near-Optimal Algorithm for Stochastic Bilevel Optimization via Double-Momentum](#)

- Prashant Khanduri, Siliang Zeng, Mingyi Hong, Hoi-To Wai, Zhaoran Wang, Zhuoran Yang
- abstract@[open-review](#): This paper proposes a new algorithm -- the \underline{S}ingle-timescale Do\underline{u}ble-momentum \underline{St}ochastic \underline{A}pprox\underline{i}matio\underline{n} (SUSTAIN) -- for tackling stochastic unconstrained bilevel optimization problems. We focus on bilevel problems where the lower level subproblem is strongly-convex and the upper level objective function is smooth. Unlike prior works which rely on \emph{two-timescale} or \emph{double loop} techniques, we design a stochastic momentum-assisted gradient estimator for both the upper and lower level updates. The latter allows us to control the error in the stochastic gradient updates due to inaccurate solution to both subproblems. If the upper objective function is smooth but possibly non-convex, we show that \{SUSTAIN\}~requires $\mathcal{O}(\epsilon^{-3/2})$ iterations (each using $\mathcal{O}(1)$ samples) to find an ϵ -stationary solution. The ϵ -stationary solution is defined as the point whose squared norm of the gradient of the outer function is less than or equal to ϵ . The total number of stochastic gradient samples required for the upper and lower level objective functions matches the best-known complexity for single-level stochastic gradient algorithms. We also analyze the case when the upper level objective function is strongly-convex.

[Generalized Jensen-Shannon Divergence Loss for Learning with Noisy Labels](#)

- Erik Englesson, Hossein Azizpour
- abstract@[open-review](#): Prior works have found it beneficial to combine provably noise-robust loss functions e.g., mean absolute error (MAE) with standard categorical loss function e.g. cross entropy (CE) to improve their learnability. Here, we propose to use Jensen-Shannon divergence as a noise-robust loss function and show that it interestingly interpolate between CE and MAE with a controllable mixing parameter. Furthermore, we make a crucial observation that CE exhibit lower consistency around noisy data points. Based on this observation, we adopt a generalized version of the Jensen-Shannon divergence for multiple distributions to encourage consistency around data points. Using this loss function, we show state-of-the-art results on both synthetic (CIFAR), and real-world (e.g., WebVision) noise with varying noise rates.

[Continual Learning via Local Module Composition](#)

- Oleksiy Ostapenko, Pau Rodriguez, Massimo Caccia, Laurent Charlin
- abstract@[open-review](#): Modularity is a compelling solution to continual learning (CL), the problem of modeling sequences of related tasks. Learning and then composing modules to solve different tasks provides an abstraction to address the principal challenges of CL including catastrophic forgetting, backward and forward transfer across tasks, and sub-linear model growth. We introduce local module composition (LMC), an approach to modular CL where each module is provided a local structural component that estimates a module's relevance to the input. Dynamic module composition is performed layer-wise based on local relevance scores. We demonstrate that agnosticity to task identities (IDs) arises from (local) structural learning that is module-specific as opposed to the task- and/or model-specific as in previous works, making LMC applicable to more CL settings compared to previous works. In addition, LMC also tracks statistics about the input distribution and adds new modules when outlier samples are detected. In the first set of experiments, LMC performs favorably compared to existing methods on the recent Continual Transfer-learning Benchmark without requiring task identities. In another study, we show that the locality of structural learning allows LMC to interpolate to related but unseen tasks (OOD), as well as to compose modular networks trained independently on different task sequences into a third modular network without any fine-tuning. Finally, in search for limitations of LMC we study it on more challenging sequences of 30 and 100 tasks, demonstrating that local module selection becomes much more challenging in presence of a large number of candidate modules. In this setting best performing LMC spawns much fewer modules compared to an oracle based baseline, however, it reaches a lower overall accuracy. The codebase is available under <https://github.com/oleksost/LMC>.

[Model-Based Episodic Memory Induces Dynamic Hybrid Controls](#)

- Hung Le, Thommen Karimpanal George, Majid Abdolshah, Truyen Tran, Svetha Venkatesh

- abstract@[open-review](#): Episodic control enables sample efficiency in reinforcement learning by recalling past experiences from an episodic memory. We propose a new model-based episodic memory of trajectories addressing current limitations of episodic control. Our memory estimates trajectory values, guiding the agent towards good policies. Built upon the memory, we construct a complementary learning model via a dynamic hybrid control unifying model-based, episodic and habitual learning into a single architecture. Experiments demonstrate that our model allows significantly faster and better learning than other strong reinforcement learning agents across a variety of environments including stochastic and non-Markovian settings.

[FedDR – Randomized Douglas-Rachford Splitting Algorithms for Nonconvex Federated Composite Optimization](#)

- Quoc Tran Dinh, Nhan H Pham, Dzung Phan, Lam Nguyen
- abstract@[open-review](#): We develop two new algorithms, called, FedDR and asyncFedDR, for solving a fundamental nonconvex composite optimization problem in federated learning. Our algorithms rely on a novel combination between a nonconvex Douglas-Rachford splitting method, randomized block-coordinate strategies, and asynchronous implementation. They can also handle convex regularizers. Unlike recent methods in the literature, e.g., FedSplit and FedPD, our algorithms update only a subset of users at each communication round, and possibly in an asynchronous manner, making them more practical. These new algorithms can handle statistical and system heterogeneity, which are the two main challenges in federated learning, while achieving the best known communication complexity. In fact, our new algorithms match the communication complexity lower bound up to a constant factor under standard assumptions. Our numerical experiments illustrate the advantages of our methods over existing algorithms on synthetic and real datasets.

[Adversarial Examples Make Strong Poisons](#)

- Liam Fowl, Micah Goldblum, Ping-yeh Chiang, Jonas Geiping, Wojciech Czaja, Tom Goldstein
- abstract@[open-review](#): The adversarial machine learning literature is largely partitioned into evasion attacks on testing data and poisoning attacks on training data. In this work, we show that adversarial examples, originally intended for attacking pre-trained models, are even more effective for data poisoning than recent methods designed specifically for poisoning. In fact, adversarial examples with labels re-assigned by the crafting network remain effective for training, suggesting that adversarial examples contain useful semantic content, just with the "wrong" labels (according to a network, but not a human). Our method, adversarial poisoning, is substantially more effective than existing poisoning methods for secure dataset release, and we release a poisoned version of ImageNet, ImageNet-P, to encourage research into the strength of this form of data obfuscation.

[Coresets for Decision Trees of Signals](#)

- Ibrahim Jubran, Ernesto Evgeniy Sanches Shayda, Ilan I Newman, Dan Feldman
- abstract@[open-review](#): A k -decision tree (t) (or k -tree) is a recursive partition of a matrix (2D-signal) into $k \geq 1$ block matrices (axis-parallel rectangles, leaves) where each rectangle is assigned a real label. Its regression or classification loss to a given matrix D of N entries (labels) is the sum of squared differences over every label in D and its assigned label by t . Given an error parameter $\varepsilon \in (0, 1)$, a (k, ε) -coreset C of D is a small summarization that provably approximates this loss to ϵ such tree, up to a multiplicative factor of $1/\varepsilon$. In particular, the optimal k -tree of C is a $(1 + \varepsilon)$ -approximation to the optimal k -tree of D . We provide the first algorithm that outputs such a (k, ε) -coreset for ϵ such matrix D . The size $|C|$ of the coreset is polynomial in $k \log(N)/\varepsilon$, and its construction takes $O(Nk)$ time. This is by forging a link between decision trees from machine learning -- to partition trees in computational geometry. Experimental results on `sklearn` and `lightGBM` show that applying our coressets on real-world data-sets boosts the computation time of random forests and their parameter tuning by up to $x10$, while keeping similar accuracy. Full open source code is provided.

[Local plasticity rules can learn deep representations using self-supervised contrastive predictions](#)

- Bernd Illing, Jean Ventura, Guillaume Bellec, Wulfram Gerstner
- abstract@[open-review](#): Learning in the brain is poorly understood and learning rules that respect biological constraints, yet yield deep hierarchical representations, are still unknown. Here, we propose a learning rule that takes inspiration from neuroscience and recent advances in self-supervised deep learning. Learning minimizes a simple layer-specific loss function and does not need to back-propagate error signals within or between layers. Instead, weight updates follow a local, Hebbian, learning rule that only depends on pre- and post-synaptic neuronal activity, predictive dendritic input and widely broadcasted modulation factors which are identical for large groups of neurons. The learning rule applies contrastive predictive learning to a causal, biological setting using saccades (i.e. rapid shifts in gaze direction). We find that networks trained with this self-supervised and local rule build deep hierarchical representations of images, speech and video.

[MobTCast: Leveraging Auxiliary Trajectory Forecasting for Human Mobility Prediction](#)

- Hao Xue, Flora Salim, Yongli Ren, Nuria Oliver
- abstract@[open-review](#): Human mobility prediction is a core functionality in many location-based services and applications. However, due to the sparsity of mobility data, it is not an easy task to predict future POIs (place-of-interests) that are going to be visited. In this paper, we propose MobTCast, a Transformer-based context-aware network for mobility prediction. Specifically, we explore the influence of four types of context in mobility prediction: temporal, semantic, social, and geographical contexts. We first design a base mobility feature extractor using the Transformer architecture, which takes both the history POI sequence and the semantic information as input. It handles both the temporal and semantic contexts. Based on the base extractor and the social connections of a user, we employ a self-attention module to model the influence of the social context. Furthermore, unlike existing methods, we introduce a location prediction branch in MobTCast as an auxiliary task to model the geographical context and predict the next location. Intuitively, the geographical distance between the location of the predicted POI and the predicted location from the auxiliary branch should be as close as possible. To reflect this relation, we design a consistency loss to further improve the POI prediction performance. In our experimental results, MobTCast outperforms other state-of-the-art next POI prediction methods. Our approach illustrates the value of including different types of context in next POI prediction.

[Early Convolutions Help Transformers See Better](#)

- Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollar, Ross Girshick
- abstract@[open-review](#): Vision transformer (ViT) models exhibit substandard optimizability. In particular, they are sensitive to the choice of optimizer (AdamW vs. SGD), optimizer hyperparameters, and training schedule length. In comparison, modern convolutional neural networks are easier to optimize. Why is this the case? In this work, we conjecture that the issue lies with the patchify stem of ViT models, which is implemented by a stride- p $p \times p$ convolution ($p = 16$ by default) applied to the input image. This large-kernel plus large-stride convolution runs counter to typical design choices of convolutional layers in neural networks. To test whether this atypical design choice causes an issue, we analyze the optimization behavior of ViT models with their original patchify stem versus a simple counterpart where we replace the ViT stem by a small number of stacked stride-two 3×3 convolutions. While the vast majority of computation in the two ViT designs is identical, we find that this small change in early visual processing results in markedly different training behavior in terms of the sensitivity to optimization settings as well as the final model accuracy. Using a convolutional stem in ViT dramatically increases optimization stability and also improves peak performance (by ~1-2% top-1 accuracy on ImageNet-1k), while maintaining flops and runtime. The improvement can be observed across the wide spectrum of model complexities (from 1G to 36G flops) and dataset scales (from ImageNet-1k to ImageNet-21k). These findings lead us to recommend using a standard, lightweight convolutional stem for ViT models in this regime as a more robust architectural choice compared to the original ViT model design.

Error Compensated Distributed SGD Can Be Accelerated

- Xun Qian, Peter Richtarik, Tong Zhang
- abstract@[open-review](#): Gradient compression is a recent and increasingly popular technique for reducing the communication cost in distributed training of large-scale machine learning models. In this work we focus on developing efficient distributed methods that can work for any compressor satisfying a certain contraction property, which includes both unbiased (after appropriate scaling) and biased compressors such as RandK and TopK. Applied naively, gradient compression introduces errors that either slow down convergence or lead to divergence. A popular technique designed to tackle this issue is error compensation/error feedback. Due to the difficulties associated with analyzing biased compressors, it is not known whether gradient compression with error compensation can be combined with acceleration. In this work, we show for the first time that error compensated gradient compression methods can be accelerated. In particular, we propose and study the error compensated loopless Katyusha method, and establish an accelerated linear convergence rate under standard assumptions. We show through numerical experiments that the proposed method converges with substantially fewer communication rounds than previous error compensated algorithms.

InfoGCL: Information-Aware Graph Contrastive Learning

- Dongkuan Xu, Wei Cheng, Dongsheng Luo, Haifeng Chen, Xiang Zhang
- abstract@[open-review](#): Various graph contrastive learning models have been proposed to improve the performance of tasks on graph datasets in recent years. While effective and prevalent, these models are usually carefully customized. In particular, despite all recent work create two contrastive views, they differ in a variety of view augmentations, architectures, and objectives. It remains an open question how to build your graph contrastive learning model from scratch for particular graph tasks and datasets. In this work, we aim to fill this gap by studying how graph information is transformed and transferred during the contrastive learning process, and proposing an information-aware graph contrastive learning framework called InfoGCL. The key to the success of the proposed framework is to follow the Information Bottleneck principle to reduce the mutual information between contrastive parts while keeping task-relevant information intact at both the levels of the individual module and the entire framework so that the information loss during graph representation learning can be minimized. We show for the first time that all recent graph contrastive learning methods can be unified by our framework. Based on theoretical and empirical analysis on benchmark graph datasets, we show that InfoGCL achieves state-of-the-art performance in the settings of both graph classification and node classification tasks.

Meta-Learning for Relative Density-Ratio Estimation

- Atsutoshi Kumagai, Tomoharu Iwata, Yasuhiro Fujiwara
- abstract@[open-review](#): The ratio of two probability densities, called a density-ratio, is a vital quantity in machine learning. In particular, a relative density-ratio, which is a bounded extension of the density-ratio, has received much attention due to its stability and has been used in various applications such as outlier detection and dataset comparison. Existing methods for (relative) density-ratio estimation (DRE) require many instances from both densities. However, sufficient instances are often unavailable in practice. In this paper, we propose a meta-learning method for relative DRE, which estimates the relative density-ratio from a few instances by using knowledge in related datasets. Specifically, given two datasets that consist of a few instances, our model extracts the datasets' information by using neural networks and uses it to obtain instance embeddings appropriate for the relative DRE. We model the relative density-ratio by a linear model on the embedded space, whose global optimum solution can be obtained as a closed-form solution. The closed-form solution enables fast and effective adaptation to a few instances, and its differentiability enables us to train our model such that the expected test error for relative DRE can be explicitly minimized after adapting to a few instances. We empirically demonstrate the effectiveness of the proposed method by using three problems: relative DRE, dataset comparison, and outlier detection.

Overcoming the curse of dimensionality with Laplacian regularization in semi-supervised learning

- Vivien Cabannes, Loucas Pillaud-Vivien, Francis Bach, Alessandro Rudi
- abstract@[open-review](#): As annotations of data can be scarce in large-scale practical problems, leveraging unlabelled examples is one of the most important aspects of machine learning. This is the aim of semi-supervised learning. To benefit from the access to unlabelled data, it is natural to diffuse smoothly knowledge of labelled data to unlabelled one. This induces to the use of Laplacian regularization. Yet, current implementations of Laplacian regularization suffer from several drawbacks, notably the well-known curse of dimensionality. In this paper, we design a new class of algorithms overcoming this issue, unveiling a large body of spectral filtering methods. Additionally, we provide a statistical analysis showing that our estimators exhibit desirable behaviors. They are implemented through (reproducing) kernel methods, for which we provide realistic computational guidelines in order to make our method usable with large amounts of data.

Unlabeled Principal Component Analysis

- Yunzhen Yao, Liangzu Peng, Manolis Tsakiris
- abstract@[open-review](#): We introduce robust principal component analysis from a data matrix in which the entries of its columns have been corrupted by permutations, termed Unlabeled Principal Component Analysis (UPCA). Using algebraic geometry, we establish that UPCA is a well-defined algebraic problem in the sense that the only matrices of minimal rank that agree with the given data are row-permutations of the ground-truth matrix, arising as the unique solutions of a polynomial system of equations. Further, we propose an efficient two-stage algorithmic pipeline for UPCA suitable for the practically relevant case where only a fraction of the data have been permuted. Stage-I employs outlier-robust PCA methods to estimate the ground-truth column-space. Equipped with the column-space, Stage-II applies recent methods for unlabeled sensing to restore the permuted data. Experiments on synthetic data, face images, educational and medical records reveal the potential of UPCA for applications such as data privatization and record linkage.

Causal-BALD: Deep Bayesian Active Learning of Outcomes to Infer Treatment-Effects from Observational Data

- Andrew Jesson, Panagiotis Tigas, Joost van Amersfoort, Andreas Kirsch, Uri Shalit, Yarin Gal
- abstract@[open-review](#): Estimating personalized treatment effects from high-dimensional observational data is essential in situations where experimental designs are infeasible, unethical, or expensive. Existing approaches rely on fitting deep models on outcomes observed for treated and control populations. However, when measuring individual outcomes is costly, as is the case of a tumor biopsy, a sample-efficient strategy for acquiring each result is required. Deep Bayesian active learning provides a framework for efficient data acquisition by selecting points with high uncertainty. However, existing methods bias training data acquisition towards regions of non-overlapping support between the treated and control populations. These are not sample-efficient because the treatment effect is not identifiable in such regions. We introduce causal, Bayesian acquisition functions grounded in information theory that bias data acquisition towards regions with overlapping support to maximize sample efficiency for learning personalized treatment effects. We demonstrate the performance of the proposed acquisition strategies on synthetic and semi-synthetic datasets IHDP and CMNIST and their extensions, which aim to simulate common dataset biases and pathologies.

Scalable Rule-Based Representation Learning for Interpretable Classification

- Zhuo Wang, Wei Zhang, Ning Liu, Jianyong Wang
- abstract@[open-review](#): Rule-based models, e.g., decision trees, are widely used in scenarios demanding high model interpretability for their transparent inner structures and good model expressivity. However, rule-based models are hard to optimize, especially on large data sets, due to their discrete

parameters and structures. Ensemble methods and fuzzy/soft rules are commonly used to improve performance, but they sacrifice the model interpretability. To obtain both good scalability and interpretability, we propose a new classifier, named Rule-based Representation Learner (RRL), that automatically learns interpretable non-fuzzy rules for data representation and classification. To train the non-differentiable RRL effectively, we project it to a continuous space and propose a novel training method, called Gradient Grafting, that can directly optimize the discrete model using gradient descent. An improved design of logical activation functions is also devised to increase the scalability of RRL and enable it to discretize the continuous features end-to-end. Exhaustive experiments on nine small and four large data sets show that RRL outperforms the competitive interpretable approaches and can be easily adjusted to obtain a trade-off between classification accuracy and model complexity for different scenarios. Our code is available at: <https://github.com/12wang3/rnl>.

Bridging Non Co-occurrence with Unlabeled In-the-wild Data for Incremental Object Detection

- NA DONG, Yongqiang Zhang, Mingli Ding, Gim Hee Lee
- abstract@[open-review](https://open-review.net): Deep networks have shown remarkable results in the task of object detection. However, their performance suffers critical drops when they are subsequently trained on novel classes without any sample from the base classes originally used to train the model. This phenomenon is known as catastrophic forgetting. Recently, several incremental learning methods are proposed to mitigate catastrophic forgetting for object detection. Despite the effectiveness, these methods require co-occurrence of the unlabeled base classes in the training data of the novel classes. This requirement is impractical in many real-world settings since the base classes do not necessarily co-occur with the novel classes. In view of this limitation, we consider a more practical setting of complete absence of co-occurrence of the base and novel classes for the object detection task. We propose the use of unlabeled in-the-wild data to bridge the non co-occurrence caused by the missing base classes during the training of additional novel classes. To this end, we introduce a blind sampling strategy based on the responses of the base-class model and pre-trained novel-class model to select a smaller relevant dataset from the large in-the-wild dataset for incremental learning. We then design a dual-teacher distillation framework to transfer the knowledge distilled from the base- and novel-class teacher models to the student model using the sampled in-the-wild data. Experimental results on the PASCAL VOC and MS COCO datasets show that our proposed method significantly outperforms other state-of-the-art class-incremental object detection methods when there is no co-occurrence between the base and novel classes during training.

A Regression Approach to Learning-Augmented Online Algorithms

- Keerti Anand, Rong Ge, Amit Kumar, Debmalya Panigrahi
- abstract@[open-review](https://open-review.net): The emerging field of learning-augmented online algorithms uses ML techniques to predict future input parameters and thereby improve the performance of online algorithms. Since these parameters are, in general, real-valued functions, a natural approach is to use regression techniques to make these predictions. We introduce this approach in this paper, and explore it in the context of a general online search framework that captures classic problems like (generalized) ski rental, bin packing, minimum makespan scheduling, etc. We show nearly tight bounds on the sample complexity of this regression problem, and extend our results to the agnostic setting. From a technical standpoint, we show that the key is to incorporate online optimization benchmarks in the design of the loss function for the regression problem, thereby diverging from the use of off-the-shelf regression tools with standard bounds on statistical error.