

Project Topic List

This course project aims to provide students with a chance to practice what they have learned in this course and deliver it as a research practice. The teaching group will provide students with some example topics to choose from. Students are also welcome to propose their own topics and link their own expertise with the knowledge taught in ECE661. Students are expected to submit their proposal and have midterm checking-in with TAs. As the final delivery, students are required to present their project during the poster session, give feedback to the Q&A session, and submit a final report. **Proposing your own project is highly encouraged.**

Please follow the project guidelines on Canvas for more information about format and grading.

All projects are expected to be done in a group of **3 students**. We expect each student to have a similar workload in the project. Each team member should be assigned balanced duty, and you should state this clearly in your proposal/report

All projects MUST adhere to ethical standards for responsible research practice and due diligence in the conduct. Please follow the Duke Community Standard. Do not use the exact same project from other courses (current or earlier) or your published work for this course.

Project Topic List

1. Input-dependent Dynamic CNN model
2. Robust and Non-robust Features
3. Improved Regularization of Convolutional Neural Networks
4. Self-Supervised / Representation Learning
5. Adversarial Patch Attack
6. Neural Network Distillation
7. Lightweight Diffusion Models
8. Learned Non-linear Quantization
9. Tiny language models
10. ***Your own proposed project (talk to TA's to finalize the details)***
11. Appendix: More projects topics from previous semesters

Project 1: Input-dependent Dynamic CNN model

Project Description

All the model compression algorithms we discuss in this course consider inducing a smaller model architecture to be used replacing the original model in the entire learning task. However, different aspects of a learning task (i.e. identifying different classes) may require totally different features. As all these features will present in the final learned model, there's no need to call all of them to deal with each individual input. This leads us to the idea of input-dependent dynamic models. In this project we will examine two lines of work in this field: dynamic kernel combination [1] and dynamic filter selection [2].

Basic Requirement

- Dynamic convolution proposed in [1] allows having more trainable parameters (convolutional kernels) in the model without significantly increasing the inference-time FLOPs. Read [1], and implement a ResNet-20 model with all the convolutional layers replaced by the dynamic convolution layer proposed.
- Train the model on CIFAR-10 dataset, start with the original ResNet-20 model, then train with 2, 3, 4 and 5 kernels in each layer respectively. Record how training and validation accuracy change as the number of kernels increases.
- Select some layers to observe the kernel attention map distribution given different inputs. Record the attention map achieved with each testing image and observe it then can be distinguished based on the class of the testing image. You may use visualization tools like T-SNE to help your observation.
- Read and implement [2] on ResNet-20, where only part of the channels is used in each layer at test time, reducing FLOPs without shrinking the model
- Report the accuracy and averaged FLOPs on CIFAR-10 validation set, visualize the gate distribution in a similar way of Figure 5 and Figure 6 in the paper

Optional Requirement (Choose the one you are interested)

- Apply white-box PGD adversarial attack to the input of [1] or [2], how is the attention map/gate selection change comparing to that of the clean input? What's the case for a transferred black-box attack?
- Perform PGD adversarial training on [1] or [2], compare performance with naïve models?
- Your own idea

References

- [1] Chen et al., “Dynamic Convolution: Attention over Convolution Kernels,” (CVPR 2020) (<https://arxiv.org/pdf/1912.03458.pdf>)
- [2] Bejnordi et al., “Batch-Shaping For Learning Conditional Channel Gated Networks,” (ICLR 2020) (<https://arxiv.org/pdf/1907.06627.pdf>)

Project 2: Robust and Non-Robust Features

Project description

It has been recently shown that both robust and non-robust features exist in common image datasets for DNN training, and both can be useful signals for standard classification [1]. DNN models tend to capture the non-robust features, which are highly correlated with the class label but are highly sensitive to small noises. Such dependency on non-robust features leads to the lack of robustness against adversarial attacks and leads to the easy transfer of attacks between models trained on the same dataset. In this project, we will use the methods proposed in [1] to disentangle the robust and non-robust features of CIFAR-10 and observe their properties.

Basic Requirement

- Train a ResNet-20 model on the CIFAR-10 dataset without adversarial training. Use the algorithm described in Appendix C.5 to generate the non-robust dataset for the CIFAR-10 training set.
- Visually observe the data under each class in the non-robust dataset. Which class do they look like? Which class does the original ResNet-20 model classify them into?
- Train a new ResNet-20 model with the non-robust dataset you got, test its validation accuracy and robustness on the original validation set.
- Train a robust ResNet-20 model on the CIFAR-10 dataset with PGD adversarial training.
- Acquire the robust feature set of the CIFAR-10 training set using your adversarially-trained model following Section 3.1 (and Appendix C.4). What do the images in the robust dataset look like?
 - Try starting from (1) a randomly sampled different image, and (2) random noise.
- Train a new ResNet-20 model with the robust dataset you got, test its validation accuracy and robustness on the original validation set.

Optional Requirement (Choose the one you are interested)

- Replace ResNet-20 model with a different ResNet model (e.g., ResNet-14, ResNet-32) by removing/adding the number of residual blocks in each stage. Repeat the experiments above. Are the robust/non-robust features related to # layers in the backbone ResNet?
- Your own idea

References

1. Ilyas et al., “Adversarial Examples Are Not Bugs, They Are Features,” (NeurIPS 2019) (<https://arxiv.org/pdf/1905.02175.pdf>)

Project 3: Improved Regularization of Convolutional Neural Networks

Project description

Convolutional neural networks (CNNs) are capable of learning powerful feature spaces. However, due to the mismatch in task complexity and model capacity, CNNs are prone to overfitting and thus require proper regularization to generalize well. When trained with standard cross entropy objectives, CNNs often struggle when the test distribution does not exactly align with the training distribution. Also, the models are often poorly calibrated, meaning that the confidence predictions they produce do not align with the true uncertainty. In this project, you will investigate and implement three techniques for improving model robustness to input noise and corruptions. You will implement the cutout regularizer [1], the mixup regularizer [2], and the self-supervised training method in [3]. After finding good configurations for these regularizers, you will compare the performance of these approaches for handling naturally occurring image noise and other corruptions.

Basic Requirement

- Train ResNet-20 [4] with cross-entropy loss (as usual) on CIFAR-10.
- Re-train the ResNet with the described cutout regularizer [1] on CIFAR-10. What are the key hyperparameters? How do these affect test accuracy?
- Re-train the ResNet with the described mixup regularizer [2] on CIFAR-10. What are the key hyperparameters? How do these affect test accuracy?
- Implement the auxiliary rotation head from [3] on a ResNet-20 model. Then train this model using the rotation head for self-supervision on CIFAR-10. What are the key hyperparameters? How much do these affect the test accuracy?
- How do these methods compare in terms of test-set accuracy? Is there a benefit to combining these techniques? What is the maximum test accuracy you can achieve?
- Test these models on a corrupted variant of the CIFAR-10 test set. Choose 5 (already implemented) corruptions from [5] to apply to each image in the test set and observe the accuracy of all models. Which performs the best? Which performs the worst?

Optional Requirement (Choose the one you are interested)

- Test robustness of the best regularized models to white-box FGSM [6] adversarial examples with various epsilon values. Do any of these techniques provide robustness to attack?
- Your own idea

References

- [1] DeVries, T., & Taylor, G. W. (2017). Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*.
- [2] Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, David Lopez-Paz: mixup: Beyond Empirical Risk Minimization. ICLR (Poster) 2018
- [3] Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, Dawn Song: Using Self-Supervised Learning Can Improve Model Robustness and Uncertainty. NeurIPS 2019
- [4] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).
- [5] <https://github.com/bethgelab/imagecorruptions>
- [6] Goodfellow, I. J., Shlens, J., & Szegedy, C. (2014). Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*.

Project 4: Self-Supervised / Representation Learning

Project description

Self-Supervised / Representation Learning is all about learning useful information from **unlabeled** data. The key is that the supervisory signal comes from the data itself, without the need for manual labeling or annotations. One method for extracting useful information from unlabeled data is to perform contrastive learning, which works by maximizing “agreement” between model output for augmented versions of the same sample. In this project, you will implement the SimCLR framework for contrastive learning of representations on CIFAR-10, then show how the model can be trained as a classifier using only a fraction of the labels.

Basic Requirement

- Implement the SimCLR [1] method for the CIFAR-10 dataset. See Appendix B.9 of the paper for detailed setup on CIFAR-10.
- Pick a reasonable batch size and number of training epochs for your hardware platform and perform a linear evaluation to reproduce the corresponding result from Figure B.7.
- Evaluate the effect of having less/limited labeled data. Follow Section 6 – Semi-supervised learning and Table 7 for inspiration of how to set this up on CIFAR10. Compare your SimCLR model to the “supervised baseline”. Is the SimCLR model better suited when there are limited amounts of labeled training data?
- Implement the RotNet [2] learning framework and compare/contrast its performance to SimCLR. Use the same two experiments described in the basic requirements for comparison.

Optional Requirement (Choose the one you are interested)

- Your own idea

References

1. Chen et al., “A Simple Framework for Contrastive Learning of Visual Representations,” (<https://arxiv.org/pdf/2002.05709.pdf>)
2. Gidaris et al., “Unsupervised Representation Learning by Predicting Image Rotations”, (<https://arxiv.org/abs/1803.07728>)

Project 5: Adversarial Patch Attack

Project description

Previously we have studied the adversarial attacks, which are realized by crafting pixel-wise additive perturbations that are constrained by L-p norm bound. These attacks, as you may notice, are not that feasible in the real-world (can you really add small noises to a physical object?). In contrast, there is another type of adversarial attacks that are designed to be “physical” – Adversarial Patch attack. As its name suggests, the idea is to create a patch (or a sticker) that can be attached to somewhere in the image to fool the classifier. In this project, you will act as an attacker and craft such adversarial patches to mislead neural networks.

Basic Requirement

- Implement the Adversarial Patch generation process described in [1] for CIFAR-10 classifiers. For simplicity you can generate rectangular patches instead of circular ones.
- Evaluate the effect of patch size. For example, generate patches with the size 3x3, 5x5, 7x7, 16x16, and measure the *untargeted* attack success rate (ASR) as a function of the patch size. Are there any visual patterns you could see in the generated patches?
- A careful attacker may want to disrupt the model in a way that whenever the patch is there the model will predict a certain target class. Repeat step 2 but this time focus on the *targeted* attack success rate. Try using each of the 10 classes as the target class. How does targeted ASR compared to untargeted ASR? Are there any visual patterns when generating targeted patches?
- Transfer the generated patches on one model to others, e.g., from ResNet to DenseNet or VGG. Measure the untargeted/targeted ASR. To what extent do the patches transfer across different models?

Optional Requirement (Choose the one you are interested)

- Your own idea

References

1. Brown et al., “Adversarial Patch,” (<https://arxiv.org/pdf/2002.05709.pdf>)

Project 6: Neural Network Distillation

Project description

In class, we have learned various compression methods to compress Neural Networks. In this project, we will investigate another way to build light weight deep networks commonly known as distillation [1][2][3]. The main idea is to let a small model, usually referred to as student model, mimic the predictions of a large model, usually referred to as teacher model. Experiments have revealed that this usually leads to better performing student models compared to solely training from ground truth labels.

Basic Requirement

- Design a small network and a large network by yourself following the instructions on section 3 in [1]. Report the FLOPS and number of parameters of both models then train them on MNIST and report their performance (accuracy) on validation dataset.
- Use the soft labels generated by the larger model to train the small network on MNIST and compare the performance with the same small model trained with ground truth.
- Try different temperature T in equation (1) in [1] and analyze the influence of this parameter on distillation.
- We refer to the dataset which the large model extracts soft label of, and the small model is trained on as transfer dataset. Obviously, the transfer dataset need not be identical to the training set of the large model. Now try to randomly omit one digit from MNIST as transfer set and repeat the distillation process (large model is still trained on full MNIST). Test the trained small model on complete validation dataset to see what the performance of this small model on the digit that it never sees during training.
- Then, try to reverse the teacher and the student model [2], which means now the student model is larger than teacher model. For this question, you could pick up two existing models by yourself as teacher and student models, such as ResNet34 and ResNet50 and experiment on Cifar10 dataset.
- Implement and Analyze KD Components [3]: First, reformulate the classical KD loss into its two parts: Target Class Knowledge Distillation (TCKD) and Non-target Class Knowledge Distillation (NCKD). Conduct an ablation study on CIFAR-100, similar to Table 1 in the paper, by training student models with only TCKD and only NCKD. Analyze your results to verify the paper's key finding that NCKD is the primary driver of performance, while TCKD alone is less effective.
- Implement and Evaluate the Full DKD Method [3]: Implement the full Decoupled Knowledge Distillation (DKD) loss, $LDKD = \alpha * TCKD + \beta * NCKD$, which is designed to overcome the limitations of classical KD. Train a student model using this method and compare its performance to a baseline student and a student trained with classical KD. The goal is to demonstrate that by decoupling the components, DKD can achieve significantly better results.

Optional Requirement (Choose the one you are interested)

- Your own idea

References

- [1] Hinton, Geoffrey, Oriol Vinyals, and Jeff Dean. "Distilling the knowledge in a neural network." arXiv preprint arXiv:1503.02531 (2015).
- [2] Yuan, Li, et al. "Revisit knowledge distillation: a teacher-free framework." (2019).
- [3] Zhao, Borui, et al. "Decoupled knowledge distillation." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022. (For DKD)

Project 7: Lightweight Diffusion Models

Project description

Generative networks are trained to model the training distribution as to generate new samples. Specifically, diffusion models are a dominant approach for high-quality generative modeling, which generate data by *reversing* a fixed noising process. In this project, you will implement a light Denoising Diffusion Probabilistic Model (DDPM) training pipeline, with DDPM sampler (full stochastic ancestral sampling) and DDIM sampler (deterministic and accelerated sampling). You will perform experiments with model design on noise schedules, parameterizations, sampler step budgets, and conditioning and editing tasks.

Basic Requirement

- **Implement a DDPM training pipeline:** First, implement the forward noising process that incrementally adds Gaussian noise to clean images. Support at least two noise schedules (e.g., linear and cosine) and provide a utility to sample a noised image at any timestep. Second, implement sinusoidal timestep embeddings, pass them through a small MLP, and inject the resulting timestep conditioning into the denoiser network. Next, implement a small denoiser network (a tiny U-Net is recommended) that takes a noisy image and the timestep conditioning and predicts the noise component. Train the model on one chosen dataset (CIFAR-10, MNIST, or Fashion-MNIST). Maintain an exponential moving average (EMA) of model weights and use the EMA weights for sampling.
- **Implement DDPM sampling:** Implement the standard stochastic reverse sampling that uses the denoiser's outputs to form a posterior estimate and then samples the next less-noisy image in the chain. Use the EMA weights for the best visual samples.
- **Implement DDIM sampling:** Implement the DDIM sampler that generates images deterministically with far fewer denoising steps than DDPM. The DDIM implementation should support variable step budgets.
- **Basic Deliverables:** (1) Training logs and loss curves showing learning progress. (2) Sampled image grids generated with: (a) DDPM (full step budget) using EMA weights, and (b) DDIM with several step budgets. (3) A documentation listing hyperparameters used and commands to run training and sampling.

Optional Requirement (Choose the one you are interested)

- **Classifier-free conditioning and guidance:** Train conditional model with label dropout and sample with different guidance scale to observe the classifier accuracy.
- **Inpainting/masked denoising:** Use model to in-paint masked regions during sampling.
- Your own idea.

References

- [1] Ho, Jonathan, Ajay Jain, and Pieter Abbeel. "Denoising diffusion probabilistic models." Advances in neural information processing systems 33 (2020): 6840-6851.
- [2] Nichol, Alexander Quinn, and Prafulla Dhariwal. "Improved denoising diffusion probabilistic models." International conference on machine learning. PMLR, 2021.
- [3] Dhariwal, Prafulla, and Alexander Nichol. "Diffusion models beat gans on image synthesis." Advances in neural information processing systems 34 (2021): 8780-8794.

Project 8: Learned Non-linear Quantization

Project description

As introduced in the lecture, a non-linear quantization scheme can reach a higher performance comparing to a linear quantization. In this project, you will be implementing one of the state-of-the-art non-linear quantization algorithms, LQ-Nets [1] (ECCV 2018), and observe its performance under different precision

Basic Requirement

- Read the LQ-Net paper and implement the quantizer and network training algorithm proposed in Section 3.3 in PyTorch.
- Use the implemented algorithm to train ResNet-20 models on CIFAR-10, with weights quantized to 3/2/1-bit. Compare your results with the results reported in Table 4.
- Apply the same quantization scheme to the activations, quantize the activations to 2-bit or 3-bit, compare the results with those reported in Table 4.
- Visualize the final weight distribution of each layer before and after quantization. How is the distribution different from linear quantization?
- Add sparsity-inducing optimization technique (e.g., soft thresholding) to the basis update step in LQ-Net algorithm, achieving mixed-precision quantization schemes by inducing zeros in the basis vector. You can only consider weight quantization for this part.
- Tune the strength of the sparsity-inducing optimization and observe the tradeoff between accuracy and average precision. How is it compared with the original LQ-Net performance?

Optional Requirement (Choose the one you are interested)

- Your own idea

References

[1] Zhang et al., “LQ-Nets: Learned Quantization for Highly Accurate and Compact Deep Neural Networks,” (<https://arxiv.org/pdf/1807.10029.pdf>)

Project 9: Tiny Language Models

Project description

Large Language Models, such as ChatGPT, are ubiquitous generative models which work well across a variety of tasks and modalities. However, their large size and computational burden lead to severe difficulties in deploying to smaller edge devices. TinyStories [1] shows that with careful construction of your dataset, even extremely small language models are capable of learning basic language semantics and generating simplistic yet consistent stories. In this project, you will explore the topic of designing and evaluating tiny language models on simple datasets. You may refer to Huggingface [2] for pre-trained checkpoints as well as the original dataset.

Basic Requirement

- Read the TinyStories [1] paper and replicate the qualitative results of Figures 9, 10, and 11. Additionally experiment with some prompts not included in the analysis of the paper. Do the models consistently produce coherent stories?
- Following the methodology in Section 4.4.1, perform quantitative analysis on both the training and validation datasets. Compare your results with those of Figures 14 – 17.
- Next, you will explore the impact of architectural choices on the learning capabilities of tiny language models. Create a few candidate architectures (~3-4) and train them on the TinyStories dataset. Repeat the evaluations performed above on these new models and compare the results between your new models. Evaluate the impact of different architectural choices.
 - Architectural choices could include modifying the number of attention heads, number of layers, matrix sizes, or you could even explore non-attention architectures such as RNNs. Make sure that the parameter count is within the range of the TinyStories models.
 - For fair comparison, keep the training parameters consistent across models.

Optional Requirement (Choose the one you are interested)

- Further simplify the TinyStories dataset. Determine the top-K most common tokens in the dataset and remove training samples which contain tokens outside of this set. Then, retrain a few TinyStories models and compare the performance. Make sure that your choice of K leaves a sizable dataset (remove no more than 50% of the original samples).
- Exploration of transfer learning capabilities for different languages.
 - Spanish: https://huggingface.co/datasets/robrenaud/multilingual_tinystories
 - French: <https://huggingface.co/datasets/MaxLSB/LeCarnet>
- Your own idea

References

- [1] Eldan, Ronen and Li, Yuanzhi, “TinyStories: How Small Can Language Models Be and Still Speak Coherent English?” (<https://arxiv.org/abs/2305.07759>)
- [2] <https://huggingface.co/roneneldan>

Appendix: More Topics from previous semesters

Given the size of the TA group, we may not be able to provide enough support for these topics. However, you are still allowed to use these topics.

Project A: Examining Lottery Ticket Hypothesis

Project Description

It is commonly believed that a larger DNN model (with more parameters or more layers) can reach a higher accuracy comparing to a smaller model under the same training scheme, thanks to the larger capacity of the model. This is the reason why people tend to prune a large DNN model rather than directly training a smaller one from scratch. However, with a specially designed initialization scheme and model architecture, a small model can possibly reach a comparable accuracy as the larger ones. Formally, this is stated as the “Lottery Ticket Hypothesis”: Dense, randomly-initialized, feed-forward networks contain subnetworks (winning tickets) that—when trained in isolation—reach test accuracy comparable to the original network in a similar number of iterations. The hypothesis as well as methods for finding the winning tickets are proposed in [1], which earn it the best paper award of ICLR 2019. You will be asked to read the paper, follow the proposed method to find the winning tickets, and examine the correctness of the Lottery Ticket Hypothesis.

Basic Requirement

- With VGG-16 or ResNet-20 model on CIFAR-10 dataset, implement the proposed algorithm for identifying winning ticket initializations, test the training performance under different sparsity
- Compare the achieved accuracy of the winning ticket with the original model, original model pruned (with iterative pruning) to the same sparsity, and randomly reinitialize the winning ticket
- Compare the difference between finding winning ticket with layer-by-layer pruning vs. global pruning, as discussed in Section 4
- Try different masking criterion other than the absolute value proposed in [1] (See [2] (NeurIPS 2019) for some ideas, you can also try saliency-based criterion introduced in the lecture or propose your own)
- Compare the results and find the best way for finding the winning tickets

Optional Requirement (Choose the one you are interested)

- Examine if winning ticket initialization can be found for PGD adversarial training
- Finding lottery ticket with Continuous Sparsification [3]
- Your own idea

References

- [1] Frankle et al., “The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks,” (<https://openreview.net/pdf?id=rJI-b3RcF7>)
- [2] Zhou et al., “Deconstructing Lottery Tickets: Zeros, Signs, and the Supermask,” (<https://arxiv.org/pdf/1905.01067.pdf>)
- [3] Savarese et al., “Winning the Lottery with Continuous Sparsification” (<https://arxiv.org/pdf/1912.04427.pdf>)

Project B: Noise-aware DNN Training

Project Description

In our course, we explored the perturbation, especially adversarial perturbation, on the input to the neural network. Meanwhile, the perturbation on the model weight is not well explored. The performance under weight perturbation actually reflects some important properties of the model. For example, the performance under L2 bounded perturbation reflects the generalizability of the model, while the performance under L infinity bounded perturbation reflects how robust the model is against quantization. This project explores sharpness aware minimization [1], which is an effective way of improving model's robustness against weight perturbation.

Basic Requirement

- Train a ResNet-20 model on CIFAR-10 dataset, evaluate its performance evaluation under weight perturbation of Gaussian noise or adversarial noise
- Start with naïve noise aware training, add Gaussian noise to the weight during training.
- Perform weight training with adversarial noise as in SAM.
- Compare performance on generalizability, quantization, and resistance against Gaussian weight perturbation of naïve SGD vs. Gaussian vs. SAM

Optional Requirement (Choose the one you are interested)

- Try SAM with multi-step weight perturbation, how is it affecting training performance?
- Try SAM with L infinity bounded perturbation. Is it more helpful on quantization?
- Your own idea

References

- [1] Foret et al., "Sharpness-Aware Minimization for Efficiently Improving Generalization,"
<https://arxiv.org/pdf/2010.01412.pdf>

Project C: Out-of-Distribution Detection

Project description

Most of the DNN classifiers we have discussed in this class are trained to be maximally accurate on the training dataset. However, is it safe to assume that during deployment the DNN will only encounter samples from the training class set? Further, what happens if a trained CIFAR-10 classifier encounters an MNIST digit? Most current networks will output a confident prediction, even though it has no real meaning. This is a realistic concern in ML operational settings and warrants research on how to distinguish “in-distribution” and “out-of-distribution” (OOD) samples at test time. In this project, you will implement several OOD detection algorithms and evaluate their performance against a variety of OOD datasets.

Basic Requirement

- Implement Baseline [1], ODIN [2], and Energy [3] methods. Your ODIN method do *not* need to include the input-preprocessing step.
- Reproduce the results in the ODIN paper using CIFAR-10 as the ID dataset and CIFAR-100, SVHN, LSUN, MNIST, Uniform-Random, Gaussian-Random as the OOD datasets. Use the AUROC and TNR@95TPR [metrics](#). Are these OOD datasets difficult to be detected?
- Try using models with different capacities (e.g., changing ResNet’s depth and width). Does the base model capacity affect results?
- There are also dedicated training methods to enable better detection. Implement Outlier Exposure [4]. Here, only use the first five classes of CIFAR-10 (denote as CIFAR-5) as ID classes, and use CIFAR-100 as the outlier dataset. After training, use the second five classes of CIFAR-10, SVHN, LSUN, MNIST, Uniform-Random, Gaussian-Random as the test OOD datasets. Compare the performance w/ and w/o OE. Does OE bring in comparable improvement for all OOD datasets?
- Evaluate OE’s performance v.s. diversity of the outlier dataset. Here we measure diversity by the number of classes contained in the outlier set. Train multiple OE models using the first 1, 5, 20, 50, 100 classes. How does the diversity affect results?

Optional Requirement (Choose the one you are interested)

- Your own idea

References

- [1] Hendrycks and Gimpel, “A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks,” (<https://arxiv.org/abs/1610.02136>)
- [2] Liang et al., “ENHANCING THE RELIABILITY OF OUT-OF-DISTRIBUTION IMAGE DETECTION IN NEURAL NETWORKS”, (<https://arxiv.org/pdf/1706.02690.pdf>)
- [3] Liu et al., “Energy-based Out-of-distribution Detection,” (<http://arxiv.org/abs/2010.03759>)
- [4] Hendrycks et al., “Deep Anomaly Detection with Outlier Exposure,” (<https://openreview.net/forum?id=HyxCxhRcY7>)

Project D: Transfer learning for domain adaptation

Project description

Semantic segmentation refers to pixel-wise classification in order to identify specific objects of interest from images capturing holistic view of the scene. Unlike typical classification tasks where one can build their own models and train from scratch, image segmentation would require rather different model architectures (such as fully convolutional network [1], FCN), and would usually need transfer learning to jump-start the model training given the more sophisticated objective [2]. This project would focus on adapting existing models pre-trained for a specific domain to another. Specifically, most existing segmentation models are pre-trained on, e.g., autonomous driving datasets, where the model segments the vehicles, pedestrians, obstacles etc. [3] from the images taken from cameras mounted around the autonomous car. On the other hand, many other domains, such as healthcare, currently have increasing needs for segmenting medical images (e.g., CT scans) in order to identify the presence of pathology, toward disease diagnoses for prognoses. In this project, you will learn how to use existing libraries to access pre-trained segmentation models, apply and fine-tune them toward a task of interest in Ophthalmology — retinal layer segmentation from optical coherence tomography (OCT) scans collected from diabetic macular edema (DME) patients. This project might be more interested in by students who are comfortable, or keen to learn, working with open-source (but non-standard) data sources and libraries online.

Basic Requirement

- Setup environment following mm-segmentation's [documentation](#), which is the library that will provide you access to many pre-trained segmentation models (to be used as the backbone), as well as the functionalities needed for transfer learning.
- Download and process the DME dataset using online [code](#), and follow the [example](#) to visualize some of the images as an initial examination of the dataset, as well as document the core statistics (e.g., how many images in total, and what are the image resolution). Then, process them in the format that can be read in by mm-segmentation, following the corresponding [documentation](#).
- Choose at least two semantic segmentation models as provided by mm-segmentation's [documentation](#), and evaluate their performance on the ophthalmic dataset, without fine-tuning.
- Based on your observation above, choose the appropriate model to be fine-tuned specifically toward achieving better segmentation performance. Unlike typical classification where the main metrics are accuracy, recall, precision etc. Segmentation outcomes are usually evaluated by metrics like intersection-of-union (IoU) score etc. to provide a holistic view of how accurate the pixel-wise classification would be. Refer to [documentation](#) and choose a few metrics that you target to improve during fine-tuning.

Optional requirement (Choose the one you are interested)

- Based on the model of your choice and fine-tuning results, do you see if any other improvements can be carried out, in order to further improve the model's performance?

Some aspects could be considered are, e.g., model architectures, other loss functions to optimize, data augmentation techniques.

- Come up with other ideas that you may want to pursue toward solving this segmentation problem.

References

- [1] Ronneberger, Olaf, Philipp Fischer, and Thomas Brox. "U-net: Convolutional networks for biomedical image segmentation." Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18. Springer International Publishing, 2015.
- [2] Kermany, Daniel S., et al. "Identifying medical diagnoses and treatable diseases by image-based deep learning." *cell* 172.5 (2018): 1122-1131.
- [3] <https://paperswithcode.com/task/semantic-segmentation>

Project E: Dynamic BlockDrop

Project description

Not all images are equally difficult to recognize. It is intuitive that some simple images can be correctly classified with a smaller model, but the difficult ones have to be classified by a larger model. Based on this motivation and the observation that ResNets behave like ensembles [1], authors of [2] proposed to learn a per-instance network routing policy that will adaptively drop certain blocks in a neural network for each image, potentially saving computational cost for the samples deemed to be ‘easy’ with minimal sacrifice of the overall accuracy. In this project we will verify the ensemble-like property of ResNets, and try to accelerate by exploiting this property.

Basic Requirement

- Train a VGG network on Cifar10 then randomly drop different layers of it and report the influence of dropping different layers on performance.
- Repeat the above experiment but this time, instead of using VGG network, try to use ResNet34 and report the influence of dropping different layers on performance (Hint: the influence of downsampling layers is different with other layers.)
- Now, we will try to drop as much as possible layers of residual networks for each image while not loss much performance. Implement a dynamic drop policy based on section 3 of [2] for ResNet34 on Cifar10.
- Test your final dynamic drop policy for ResNet34 on Cifar10. Then compare randomly drop policy and dynamic drop policy based on dropping ratio (how many blocks you drop) and performance (accuracy on validation dataset.)
- Based on your results, could you provide several reasons for explaining why residual networks are better than previous networks, such as VGG?

Optional Requirement (Choose the one you are interested)

- Your own idea

References

- [1] Veit, Andreas, Michael J. Wilber, and Serge Belongie. "Residual networks behave like ensembles of relatively shallow networks." *Advances in neural information processing systems* 29 (2016): 550-558.
- [2] Wu, Zuxuan, et al. "Blockdrop: Dynamic inference paths in residual networks." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2018.

Project F: Deepfake Generation

Project description

As deep learning becomes more and more powerful, deep learning based deepfake generation methods can produce more and more realistic-looking deepfakes. Fake face images can be generated by either face synthesizing or face manipulation. Face synthesizing creates new faces that do not belong to any person in the world, while face manipulation manipulates the face images of existing individuals. Generally speaking, face manipulation may result in more severe consequences than face synthesizing, because it may target at individuals who have high social impact or make illegal use of other people's faces. Therefore, in this project, we will focus on face manipulation based Deepfake generation and discuss the potential ways to defense against it.

Basic Requirement

- Prepare datasets for Deepfake generation. Here's some typical datasets for you to consider: FaceForensics++[1], DFDC[2]. Feel free to search for other datasets. You can also use your own face images and your friends' (should be fun ^_^).
- Face manipulation can be further categorized into face replacement (face swap) and face reenactment. In face replacement, a target face is replaced as a source face and the generated fake face claims to be the source. In face reenactment, the expression, hair color, and/or other properties of the target face are changed as those of a source face, but the identity of the target face is still preserved. Implement both face replacement and face reenactment. Show the fake faces you generated and comment on the results.
- Existing deepfake detection mainly focuses on passive detection. Train a classifier using the true faces and fake faces you generated to detect whether a given face is deepfake. Try many different classification models, compare their performances and select the best one. Select a proper metrics to measure the performance of the classifier and report the results.
- Use another dataset to train your deepfake models and generate fake faces. Then use the classifier you trained before to detect true/fake faces (Do not re-train the classifier on new fake faces). Can this classifier still correctly identify the fake faces?
- Do you think it's a good way to defense against deepfake? Discuss different methods to defense against the deepfake or reduce the harm of deepfake.

Optional Requirement (Choose the one you are interested)

- Your own idea

References:

- [1] [FaceForensics++: Learning to Detect Manipulated Facial Images](#)
- [2] [The DeepFake Detection Challenge \(DFDC\) Dataset](#)