

Robert Tjarko Lange - Handwritten Notes

08/2019

MATH+ / BMS Summer School "Mathematics of Deep Learning" 2019 Schedule

Time	Monday 19.8.	Tuesday 20.8.	Wednesday 21.8.	Thursday 22.8.	Friday 23.8.
8.00	Registration				
8.45 - 9.00	Welcome				
9.00 - 10.30	Kutyniok 1	Fleuret 2	Vidal 3		6. Bergman Spokoiny
10.30 - 11.00	Coffee	Coffee	Coffee		PDE / AV
11.00 - 12.30	Fleuret 1	Vidal 2	Haber 2	→ 6 hours of Habers ODE	TGIF Meeting
12.30 - 14.00	Lunch	Lunch	Lunch	Lunch	Lunch
14.00 - 15.30	Vidal 1	Haber 1	Fleuret 3	Haber 3	
15.30 - 16.00	Coffee	Coffee	Coffee	OAE / PPE PoV	Coffee
16.00 - 17.30	Poster Session	Tensor Flow (Basics)	Tensor Flow (Basics)	Kutyniok 2	

Inverse Problems

Time	Monday 26.8.	Tuesday 27.8.	Wednesday 28.8.	Thursday 29.8.	Friday 30.8.
9.00 - 10.30	Müller	Jenssen 2	Krause 1	Cohen 3	Leimkuhler 3
10.30 - 11.00	Coffee	Coffee	Coffee	Coffee	Coffee
11.00 - 12.30	Jenssen 1	Debate: New Horizons	Cohen 2	Schütte 2	Kutyniok 3
12.30 - 14.00	Lunch	Lunch	Lunch	Lunch	Closing
14.00 - 15.30	Leimkuhler 1	Schütte 1	Krause 2	Jenssen 3	
15.30 - 16.00	Coffee	Coffee	Coffee	Coffee	
16.00 - 17.30	Cohen 1	Practice Session	Leimkuhler 2	Practice Session	

to start on Thursday 2nd of August → before Lab Roberts!

3rd Lecture Haber (UBC) - Using PDE to design stable architectures

stepsize

- DNN as ODEs : $y_{j+1} = y_j + \Omega K^{(1)}_j$ or $(K^{(1)}_j y_j + b_j)$ \rightarrow time steps \Leftrightarrow layers!
 $y = K^{(2)}(t) \circ (K^{(1)}(t) Y + b(t))$
- Regularization & early stopping \Rightarrow badly-posed forward problem!
- Other problem \rightarrow choose $K^{(1)}, K^{(2)}$ appropriately \rightarrow Hamiltonian!
- Optimize \Leftrightarrow discretize \Leftrightarrow local order? ! SO vs. ODE

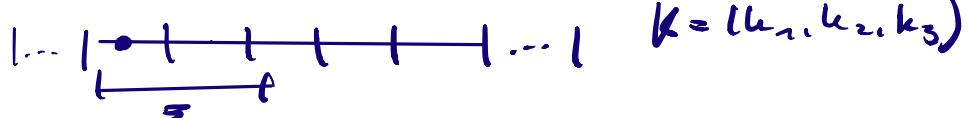
↳ Classical: First discretize then optimize

↳ Why otherwise? Henry Cudworth!

- Meshless methods \rightarrow continuous spaces \rightarrow like & we pre-^{→ better} process!
- Point on the collocation points \rightarrow branches usually hidden!
- ↳ Does Colvad have a bias? ; the batch norm seems tractable?
- ↳ Example: large batch \Rightarrow allows for larger step-size
- ↳ $\partial t \rightarrow$ stepsize! , y_j time steps \Leftrightarrow "layers" fixed number of elements
- ↳ Problem! No degrees of freedom locally \rightarrow Cost! Sol example!
- ⇒ Under how to subdivision / number of features in ODEs!**
- Forget as much as possible! No dear best try with practical example
- DL like every inverse problem \rightarrow look at intermediate steps
- Sanity check: Make sure that initial loss is close to random!
- Small batches initially \rightarrow slower as time \Rightarrow problem: Of parallelization of a batch
- Monitor how gradient behaves: $g \leftarrow \| \nabla \ell \| \rightarrow$ look at change!
- CNNs as PDEs: $K^{(1)}, K^{(2)}$: Conv kernels

\rightarrow Conv in 1D:

$$\begin{pmatrix} k_2 k_3 & 0 & k_1 \\ 0 & k_3 k_2 & 0 \\ 0 & 0 & k_1 k_2 \end{pmatrix}$$



\rightarrow degree of basis! \Rightarrow approx by depth

→ PDE → Space + time!

→ Parabolic PDEs: $y_t = -K(t)^T \sigma [K(t)y + b(t)]$

↳ Reduced form: $y_{j+1} = y_j + h K_j^{(1)} \sigma (K_j^{(1)} y_j + b_j)$

with constraint of $K_j^{(1)} = -(K_j^{(1)})^T \Rightarrow$ weight

↳ Van-Linear heat equation relationship! \Rightarrow provides us with an algorithm

→ Weight decay \Leftrightarrow Tikhonov regularization

• Parabolic model / net: backward!

$$y_{jt} = -K(t)^T \sigma [K(t)y + b(t)]$$

$$y_{j+1} = 2y_j - y_{j-1} - h^2 K_j^T \sigma (K_j y + b_j)$$

→ Reversible: forward + backward in time!

→ less memory required for backprop \Rightarrow can reconstruct activations

backwards from y_{j+1} and y_j \rightarrow don't need to store all activations

• Full West and semi-implied networks! \Rightarrow Diffusion - Backprop process

• PDE changes should be based on problem at hand!

• Batchsize not related to h \rightarrow controls what net can express!

\rightarrow more freedom vs. optimizability!

• Neural ODE paper relationship \Rightarrow First opt. Then discretize

↳ focus that you can backprop interpret \Rightarrow but can't!

↳ Fundamental problem in control theory!

2nd Lecture Kutyniok (TU Berlin) - DL meets Inverse Problems

- Talk outline: Between model-based and data-driven approach
- Inverse problem = Recover original data from transformed version
 - Examples: Fill in missing data, Feature extraction & Regularization
 - Formulation:

$$K: V \rightarrow Y \text{ with } Kx = y \quad \forall x, y \in X, Y$$

→ Hardened conditions \Rightarrow Hardened

(1) Existence, (2) Uniqueness: $x \in X$ is unique (3) Stability

→ Regularization before uniqueness! \rightarrow PROR! Generalization

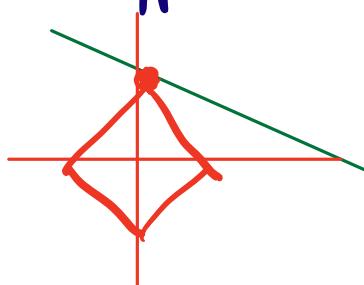
- Tikhonov Regularization: An approx solution $x^\alpha \in X$, $\alpha > 0$

$$\min_x J_\alpha(x) := \|Kx - y\|^2 + \underbrace{\alpha \|x\|^2}_{\text{data fidelity term}} \quad \underbrace{x \in X}_{\text{reg. term}}$$

→ Generalization \Rightarrow why worse? $\rightarrow \alpha \cdot \underline{R(x)}$

↳ more continuous dependence on data

↳ Different examples: $L^2 \Rightarrow$ smallest energy



$L^1 \Rightarrow$ piecewise constant functions

$2^1 \Rightarrow$ Sparsity inducing

all solutions to inverse problem

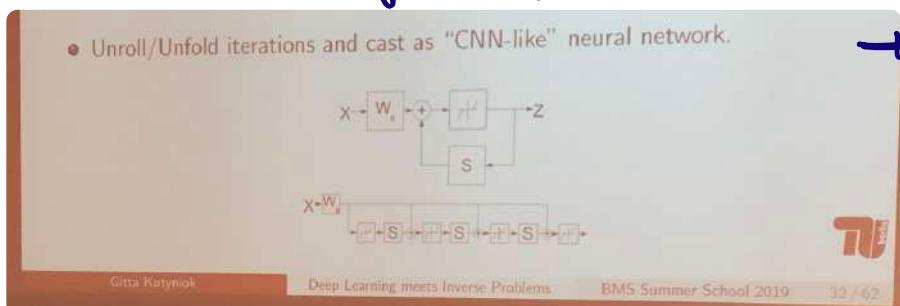
- Sparsity:
 - Sparse signals vs. compressible signals
 - ↳ number non-zero coefficients of exact
 - ↳ rapid decay in size of coefficients \Rightarrow per few

- Compressed sensing:
 - $y = Ax \rightarrow A \gg n \Rightarrow$ Sparse solution after transform

↳ Compressive sensing in single step \Rightarrow solve via regularized!

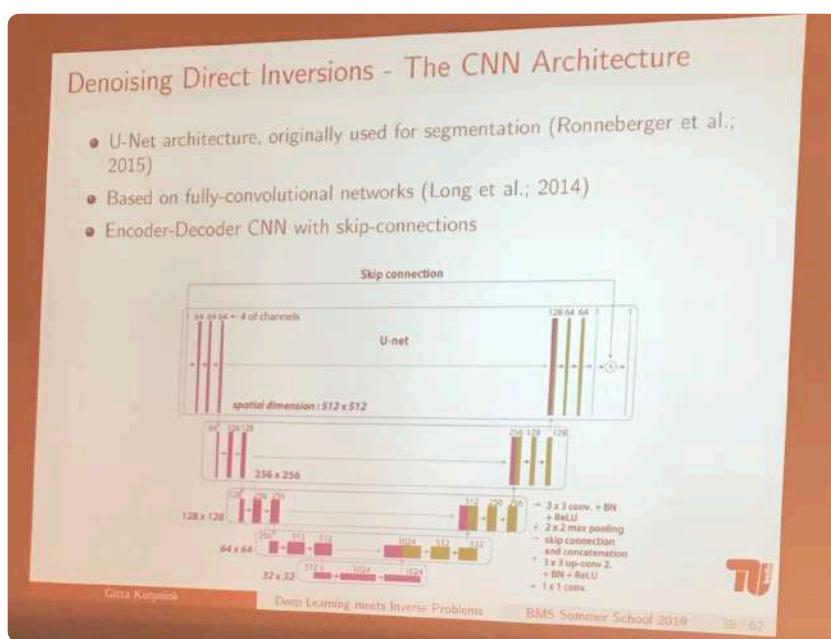
- Frame \Leftrightarrow generalization of orthonormal basis
 $\rightarrow (\psi_d)_d ; \exists 0 < t \leq B < \infty$ with $t\|x\|^2 \leq \sum_d |\langle x, \psi_d \rangle|^2 \leq B\|x\|^2$ learn!
- Redundancy: Robust to loss of dimensions/basis! \Rightarrow can be transmitted well
- Find such a system \Rightarrow dictionary (representation syst.) learning algorithms, for noise reduction! wavelets, shearlets L-SVD!

- First iteration: $x_{n+1} = \sum_k [x_n - \frac{1}{2} \cdot w_k^T (w_k x_n - y)]$
↳ soft-threshold
 for solving $\min_x \|w_k x - y\|_2^2 + \gamma \|x\|_1$



→ Formulate in network structure!
 → Leber et al (2019)

- Denoising w. additive Gaussian Noise \rightarrow Solved (Eckel et al., 16')
- Denoising Direct Inverses (Ye et al.; 2018): first reconstruct and then denoise reconstructed image
↳ Classical architecture: U-net



• Procedure to solve Filchner:
Poisson - Radon
↳ also related to denoising
⇒ soft-threshold

- Often times nowadays pre-trained denoising CNN is taken and plugged into denoising parts of applications!
- Alternative approach via generative models \Rightarrow fNNs // Applicable to CT scans

Lecture - Leonid Berlyand (Penn State) - 'PDE techniques in DL'

Outline of Talk

- ① Classification problem & DNN-based algorithms: Definitions and Issues
 - (a) Introducing DNN-based classifiers and their training
 - (b) Issues of training: accuracy, convergence, stability & robustness
- ② PDE (kinetic theory) & variational (min-max) techniques for DNN-based algorithms:
 - (a) Convergence result [LB, Jabin: Comptes Rendus - Math, 2018]
 - (b) Stability result [LB, Jabin: in progress, 2019]
- ③ Current work:
 - (a) Improving conditions for stability [LB, Alex Safsten (PSU), & Jabin]
 - (b) Numerics & analysis for accuracy [LB, Robby Creese (PSU) & Jabin]
 - (c) Low dimensional manifold learning [LB, M. Potomkin (PSU), Nikita Puchkin (Moscow), V. Spokoiny (WIAS)]

- Stability \rightarrow Discretize \Rightarrow Numerics, Computation from ANALYSIS point to disentangle contributions
- $\phi: \mathbb{R}^n \rightarrow [0, 1]^m \Rightarrow$ Non-exact $\phi(s, \omega)$ to approximate classifier \Rightarrow Opt. + priors!
- $\langle L(x, s) \rangle_T$

- Rate of converge: $\|x^{n+1} - x^n\| < f(n)$
- Stability \leftrightarrow dynamics of classification accuracy
- Universal approx. Th.: Single Hidden Layer Net w. inf. width can approximate smooth func. \Rightarrow Only EXISTENCE!

$$\bar{L}(x) = - \sum_{s \in T} \frac{r(s)}{\text{true weight of class}} \log p_{\text{class}}(x, s) = \sum_{s \in T} r(s) L(x, s)$$

true weight of class $\Rightarrow [0, 0, 1, \dots, 0]$

- Accuracy + Loss! \rightarrow Min-max problem
- PDE perspective via SGD
 - \rightarrow Stochasticity does not allow you to take continuous time limit
 - \rightarrow Requires kinetic PDE approach \rightarrow entropy - entropy dissipation

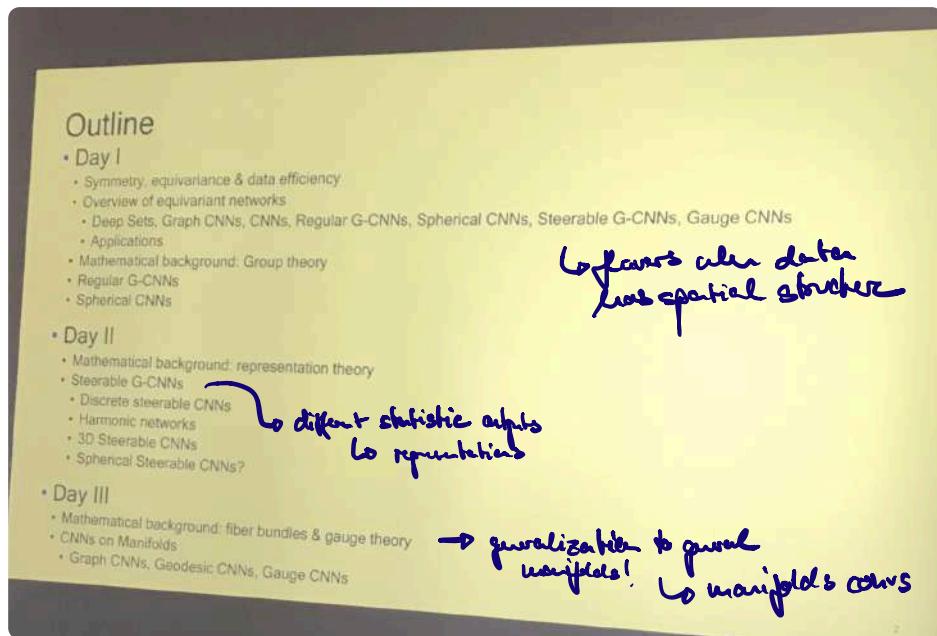
$$\frac{dE}{dt} = -D(t) \quad \begin{aligned} \rightarrow E: \text{entropy} \rightarrow E(S) := \bar{L}(\tilde{x}(S)) \\ \rightarrow D: \text{dissipation} \rightarrow D(S) := |\nabla_x \bar{L}(\tilde{x})|^2 \end{aligned}$$

To Orouwall bound: $E(S) = \bar{L}(S) \leq C_f \rightarrow 0$ as $t \rightarrow \infty$

\rightarrow Fokker-Planck Equations: Allows to formulate stat. prob. w.r.t. pdf form $f_{\text{true}}(x) = -\Lambda_T^{-1} \tau_{\text{true}}(x)$

1st lecture - Taco Cohen - Equivariant Networks

- Prior knowledge of structure of symmetry of problem \Rightarrow data efficiency



• Symmetry = Transform class
not single object

\hookrightarrow ML: * Prob. distr. \rightarrow Sampling in multi-modal dist.

* Label fcts. $X \mapsto Y$

* Permutation fcts.
 \rightarrow permutation of neurons and their weights

* Causality = inverse to shifts in digits.

\hookrightarrow Symmetry as inductive bias \rightarrow constraints define laws of physics

\hookrightarrow Statistics: De Finetti \Rightarrow Exchangeability \leftrightarrow Conditional Independence

• Transfo Group : Set of transfo \rightarrow identity, composed associatively, invertible

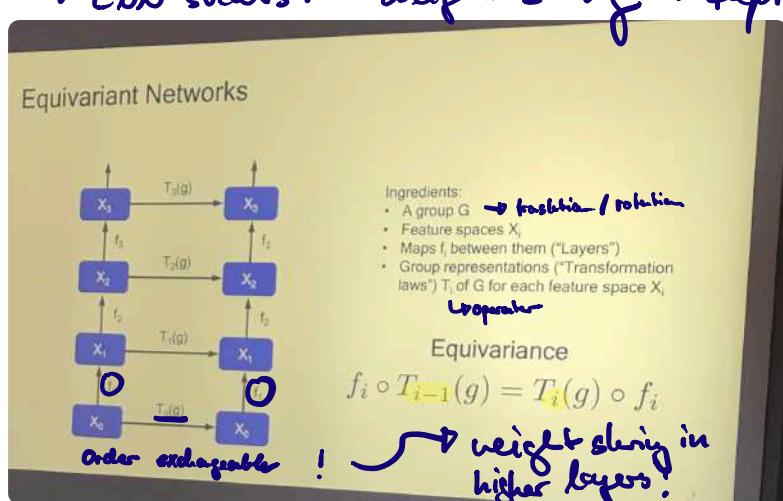
\hookrightarrow Rotations (filters, feature maps), Euclidean motions, Projections, Diffeomorphisms (continuous, discrete)

• Picasso Problem: Invariance not enough \rightarrow Individual parts are not enough to tell you the identity of an object

\rightarrow CNN success: weight sharing + depth + Translation Equivariance

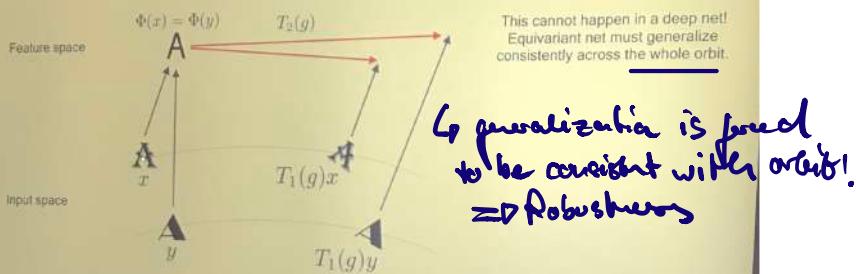
$\rightarrow f_1 \circ f_2$ equiv. $\Rightarrow f_1 \circ f_2$ equiv.

\hookrightarrow To build an Equiv Net:
make sure that all layers respect symmetries / are equivariant



Equivariance as Symmetry-consistent generalization

$$\text{Equivariance: } \Phi(T_{i-1}(g)x) = T_i(g)\Phi(x)$$



- Data augmentation
→ does not guarantee equivariance!!
→ Combinatorial explosion in number of symmetries

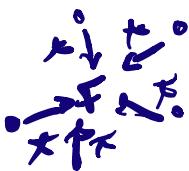
- Residual connections → Equivariance also works again!

Equivariance to Symmetry Transformations

Examples

Data	Images, Audio, ...	Graphs	Sets	Signals on homogeneous space	Signals on manifold
Symmetries	Translations; Rotations	Permutations	Permutations	Group G	Gauge group
Architecture	CNNs; Group CNNs	Graph NNs	Deep Sets	Group CNNs	Gauge CNNs

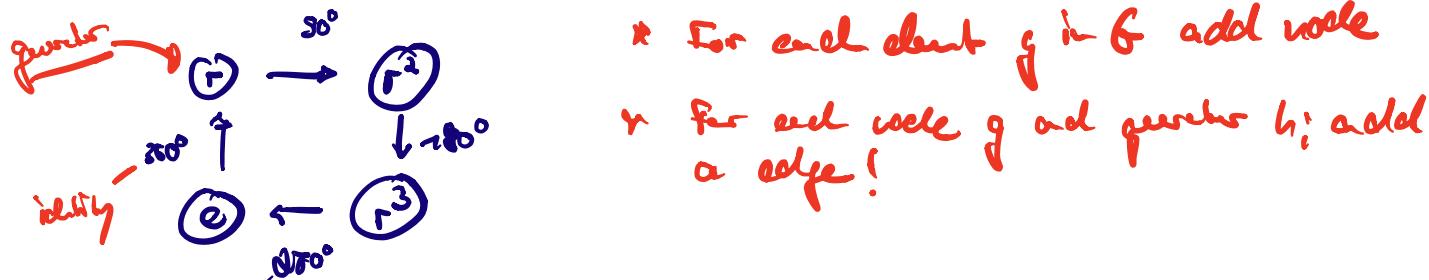
- Deep sets → Zalik et al 17'
→ ordering does not matter → Removable
→ Symmetry is trivially independent of ordering
- Graph Convolution via Message Passing



- CNNs → filters only use \downarrow & \leftrightarrow
→ Non-rotational equivariant
→ Graph - Convolution \Rightarrow 4- 90° -chains

- 3D R-CNNs → 3d Cables → it handles symmetries implicit!
- Scale Equivariance! → generally non-invertible → semi-group env.

- CNNs on Manifolds
 - Local neighborhoods via tangent space
 - ambiguity in choice of filters \Rightarrow gauge symmetry \rightarrow ~~Mark Pol~~
- Group Theory Basics
 - Cayley Diagrams \Rightarrow group representation as a graph



- Cosets: For $g \in G, h \in H: g \sim g' \Leftrightarrow \exists l \in H: g' = gh$
↳ reflexive, symmetric & transitive
↳ Equivalence classes = Cosets \Rightarrow Visualize in Cayley!
↳ Quotient
- Group action: Maps $(g, x) \mapsto gx$
- Homogeneous space = Quotient space w/o privileged origin/cosets

CNN	G	H	G/H
Conventional CNN	\mathbb{Z}^2 (planar translation)	$\{e\}$ (trivial)	\mathbb{Z}^2 (2d pixel grid)
P4-CNN	P4 (discrete roto-translation)	C4 (4-fold rotation)	\mathbb{Z}^2 (2d pixel grid)
SE(2)-CNN	SE(2) (continuous roto-translation)	SO(2) (continuous rotation)	\mathbb{R}^2 (continuous plane)
Spherical CNN	SO(3) (continuous 3D rotation)	SO(2) (continuous 2D rotation)	S^2 (sphere)
DNA CNN	$\mathbb{Z} \times \mathbb{Z}_2$ (translation + inverse-complement)	\mathbb{Z}_2 (reverse-complement)	\mathbb{Z} (1d grid)

- Regular G-CNN
- Group Conv: Output feature map ab $g \in G$: Transform filter by g and compute inner product with input feature maps
- feature map: $f: X \rightarrow \mathbb{R}^C$ \Rightarrow think of as filter / channels stored for one particular location!
- Filter bank / Kernel: $\mathcal{F}: X \rightarrow \mathbb{R}^{K \times C} \Rightarrow$ Support small: 3×3 e.g.

→ Prayes of feature maps & filters : Group actions on feature maps/filter

$$[L_g f](x) = f(g^{-1}x)$$

→ Group Commutative: $[\phi * \tau](g) = \langle L_g \phi, \tau \rangle$

Feb. on
group

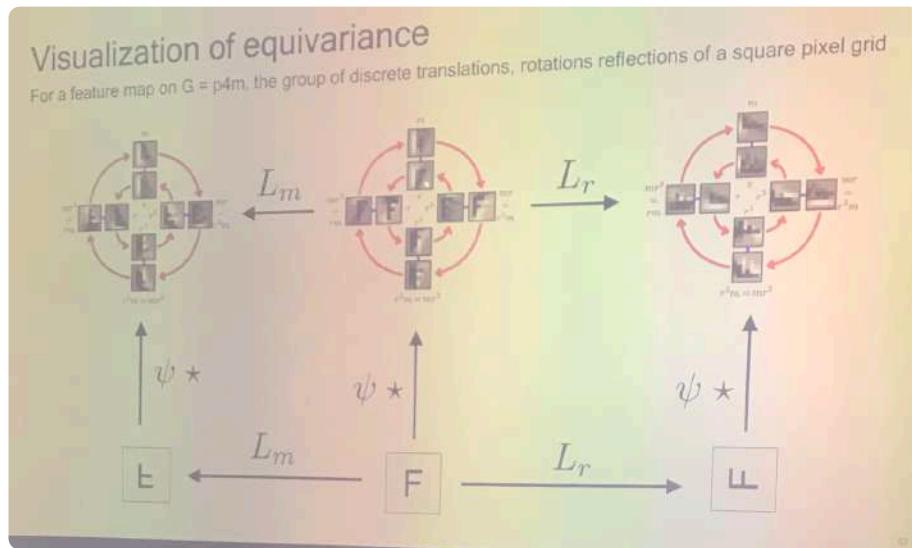
trans
for later

intro
pet.

$$= \sum_{x \in X} \sum_{c=1}^C p_c(g^{-1}x) f_c(x)$$

Visualization of equivariance

For a feature map on $G = p4m$, the group of discrete translations, rotations reflections of a square pixel grid



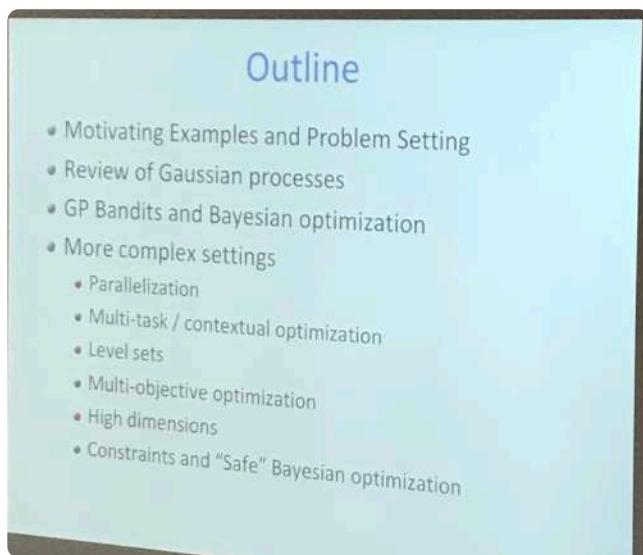
- Implementable
 - Filter bank just bigger!
 - ↳ index for rotebox
 - ↳ filter bank step!

↗  2x2 Filter
⇒ Rotebox!

- Non-linearity $[C_r f](g) = [r \circ f](g) = r(f(g))$
 - Proper / Radonity \rightarrow loss of equivariance!
 - Trainability \rightarrow Gradient on propagable graph all without filter banks!

1st Lecture - Andreas Krause - Bayesian Optimization

- Protein Filter Example → Collaboration w. Chemistry nobel prize winner Horst Hahn
 - ↳ SEQUENTIAL EXPERIMENTATION! w. MOLSY feedback
 - ↳ E-E Dilemma → NEED TO USE REGULARITY



- h-armed stochastic bandits
 - allocation of T tokens
 - payout over $i \Rightarrow \text{mean } f_i$
 - ↳ D-armed generalization!
 - Can't try all options one
 - Need to generalize/predict
- structured bandits : Linear, Lipschitz, combinatorial, etc.
- Original BO paper → Hockeys 75'
 - ↳ less theoretical background compared to bandits
 - ↳ real problem is very similar!

Cumulative Regret: $R_T = \sum_{t=1}^T (\max_x f(x) - f(x_t))$

↳ smaller if $R_{T+1} \rightarrow 0$

Simple Regret: $S_T = \min_{t \in \{1, \dots, T\}} (\max_x f(x) - f(x_t))$

↳ Min s.s. b.v.: $S_T \leq R_T/T$

⇒ Q: Does the max always have to be linear?

- Gaussian Process \Rightarrow 0D-dim multivar. Gaussian \Rightarrow over fct.
 \hookrightarrow finite angulars are multivar. Gaussians space

$\mu: X \rightarrow \mathbb{R}$, $K: X \times X \rightarrow \mathbb{R} \Rightarrow$ PRACTICALLY.

$\hookrightarrow N(\mu_x, \Sigma_{xx}) \Rightarrow K$: kernel, μ : new fct.

\hookrightarrow Cov./Kernel characterizes smoothness abs.!

- positive definiteness requirement of kernels CORRELATIONS

- $K(x, x') = \exp\left(-\frac{(x-x')^2}{l^2}\right)$ slight scale \rightarrow wigglier!

- linear kernel $\Rightarrow K(x, x') = \begin{pmatrix} x & x' \end{pmatrix}^\top$ \Rightarrow Bayesian Linear Regression
 $\hookrightarrow \Phi(x)^\top \Phi(x') \rightarrow$ Learn!

SOUTHEAD
EXP.

Making predictions with GPs

• Suppose $P(f) = GP(f; \mu, k)$

and we observe $y_i = f(x_i) + \epsilon_i \quad \epsilon_i \sim N(0, \sigma^2)$

• Then the posterior is also a GP: $A = \{x_1, \dots, x_n\}$

$$P(f | x_1, \dots, x_n, y_1, \dots, y_n) = GP(f; \mu', k')$$

$$\mu'(x) = \mu(x) + \Sigma_{x,A} (\Sigma_{AA} + \sigma^2 I)^{-1} (y_A - \mu_A)$$

$$k'(x, x') = k(x, x') - \Sigma_{x,A} (\Sigma_{AA} + \sigma^2 I)^{-1} \Sigma_{A,x'}$$

Kernel Characterization

- Prior knowledge

- Max ML. \Rightarrow Empirical Bayes

\hookrightarrow Interpretation hyperp.

- Online adaptation!

\hookrightarrow Form of Bayesian theory

\Rightarrow have to use non-conjugate prior!

\hookrightarrow Bandits \Rightarrow surrogate model per fct. of f \rightarrow How can we cast shade the fct./GP when we know what the max. is?!

• Uncertainty quantification?!

\rightarrow Entropy: $H(x) = E[-\log p(x)] \Rightarrow$ Max info gain/entropy reduction

$$\hookrightarrow f(s) = H(f) - H(f|y_s) = \frac{1}{2} \log |1 + e^{-2\sum s_i}|$$

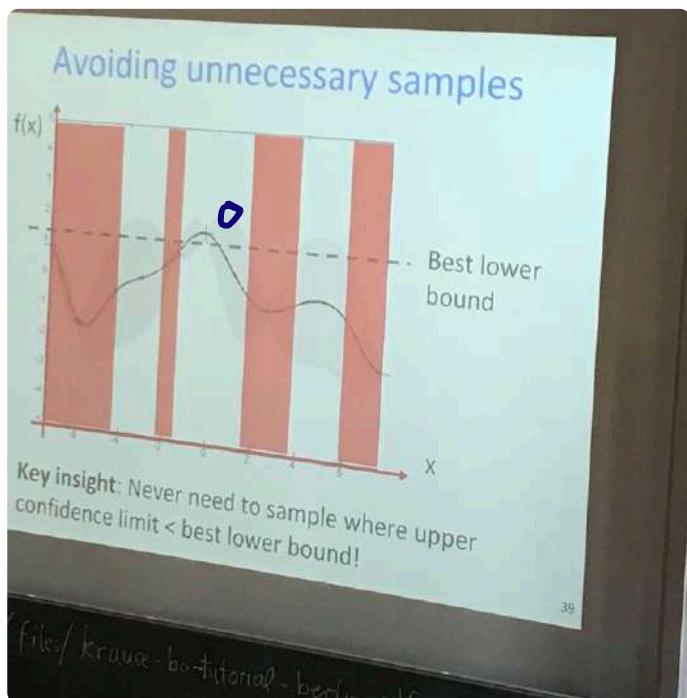
↙ ↓ ↘
 entropy before eval of y_s entropy after eval of y_s at selected point!

\hookrightarrow UP level to optimize

→ Greedy: $f(s)$ is monotone submodular $\Rightarrow s = \max_j \sum s_j$

$$\forall x \in D, \forall A \subseteq D \subseteq D: f(A \cup \{x\}) - f(A) \geq f(D \cup \{x\}) - f(D)$$

↪ Only explore \rightarrow less about problem but waste a lot of samples in part of space where problem is exp. to be small \rightarrow OPTIMIZATION $\rightarrow E \text{ vs. } E$



- UCB: Take best mean value at the current state μ^{best}

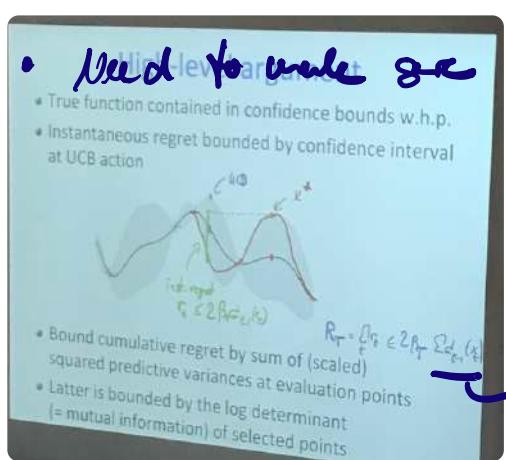
$$\hookrightarrow x_t = \arg \max_{x \in D} \mu_{b+1}(x) + \beta \frac{\sigma_{b+1}(x)}{t}$$

↳ increase our tree \rightarrow fine-grained since variance decreases over time

↳ Max vs. Min! \Rightarrow Our objective is to max profit!

↪ Can derive converge!

- Need to make sure that the feasible lies within action bands!



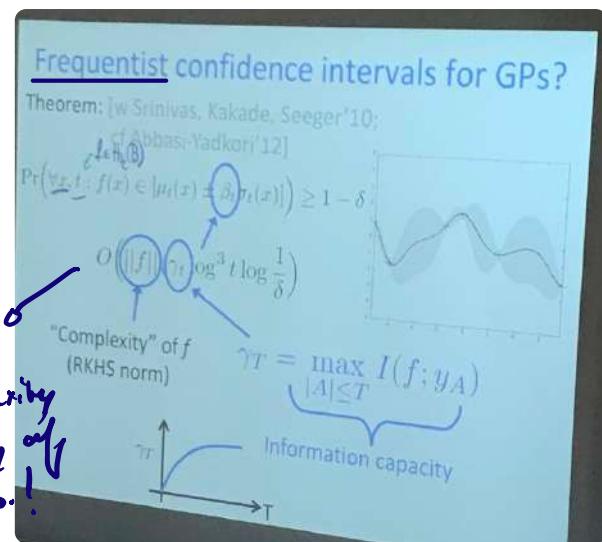
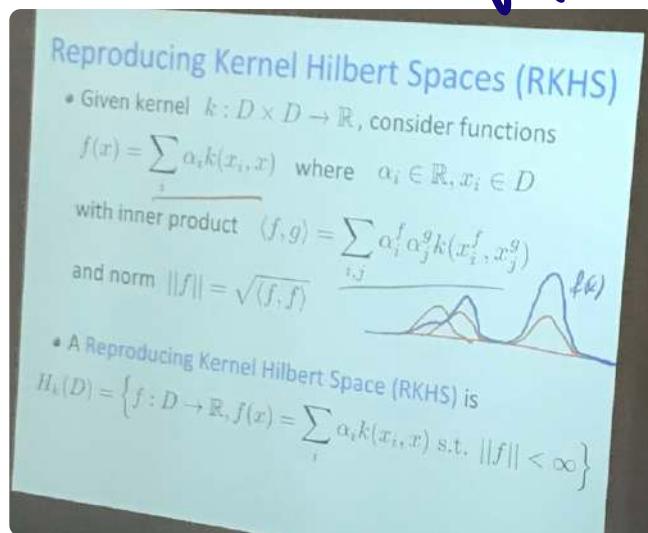
→ $f_T \Rightarrow$ estimate with confidence

$$\text{where } f_T = \max F_T$$

Variance contracts over time!

↪ based on avg. regret $O(\sqrt{T})$

- Need 2 crucial steps: 1 True function with bands 2 max f(x) b/w
- Robustness: Family of functions \rightarrow Want $f(x) \in [f_{\text{low}}(x) \pm \delta_{\text{low}}(x)]$



\rightarrow Generally speaking: Non-convex opt. of app. fct.!

\hookrightarrow Lipschitz - optimization \hookrightarrow Gradient ascent!



2nd Lecture - Andreas Krause - Extensions

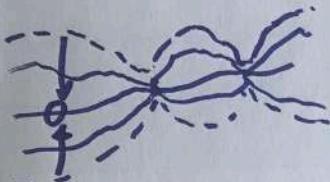
Andreas Krause - BO/GP Bandits - Lecture 2

- Parallelization = Delayed feedback \rightarrow Batch Selection
 \rightarrow internal noise \rightarrow also value in doing same exp. multiple times

$$x_t = \underset{x \in D}{\operatorname{argmax}} \left[\mu_{t-1}(x) + \beta_t \sigma_{t-1}(x) \right]$$

anticipate variance reduction!
"hallucinate"

\hookrightarrow idea: inflate confidence bands such that no matter what happens during experiments \rightarrow they are still recursive!



\hookrightarrow Then hallucinate for variance reduction
 at 1st point in batch, afterwards:
 Select next point, hallucinate var. red., ...

\hookrightarrow Google魏氏 AutoML

\rightarrow Other: Schonlau 97!: Multi-Point EI
 Azimi et al 10!: Simulation Matching

- Multi-Task / Transfer BO \rightarrow Contextual GP bandits

$$\rightarrow \text{context var. } \Xi_t \in \Xi \rightarrow y_t = f(x_t, \Xi_t) + \epsilon_t$$

$\hookrightarrow r_t = \sup_k f(x, \Xi_t) - f(x_t, \Xi_t)$

\hookrightarrow borders: map context to action \rightarrow between bandit and RL!
 \rightarrow not model transitions / sequence of contexts \rightarrow arbitrary

$$y_t = \underset{x \in D}{\operatorname{argmax}} \mu_{t-1}(x, \Xi_t) + \beta_t \sigma_{t-1}(x, \Xi_t)$$

\hookrightarrow smoothness along context-action pairs!

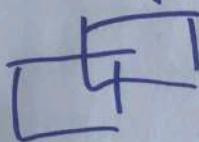
$$k((x, \Xi), (x', \Xi')) = k(x, x') \cdot k(\Xi, \Xi')$$

or addition

\rightarrow Sublinear Contextual regret bounds!

Other methods?

- Multiple Objectives \rightarrow Pareto Optimality Consideration
 - \rightarrow Predict how pareto front looks like from few evals
 - \hookrightarrow As for different objectives! \rightarrow Different GPs!
 - \hookrightarrow Allow to rule out several evaluable points!
 - \rightarrow allow for study in pareto classification



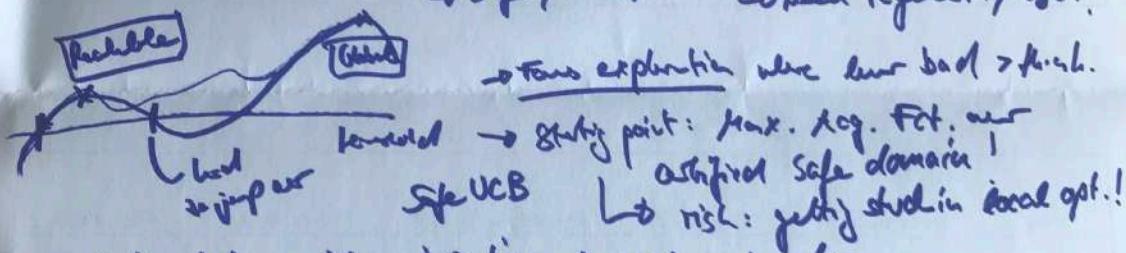
\rightarrow Sample from regions most informative!

PARETO

Discrete Vars.
= Graph kernels
= Learn representativity

- Safe Exploration + Constraints \Rightarrow PBS / Stabilization
 - \rightarrow Sim & Real

$$\max_{\theta} f(\theta) \text{ s.t. } g(\theta) \geq 0 \quad \begin{array}{l} \rightarrow \text{Max. Value fn.} \\ \rightarrow \text{value function} \\ \hookrightarrow \text{safety model} \quad \hookrightarrow \text{constraints} \\ \hookrightarrow \text{Need regularity ass.} \end{array}$$



- \rightarrow Need to provide incentive to explore boundary
 - \hookrightarrow pick aters that work $g \rightarrow$ same type of critics

\rightarrow Need to transform time / physical costs into bits!

$$\hat{f}(x) = f(x) + \underbrace{\delta f(x)}_{\text{bits induced by sim and not real!}}$$

\hookrightarrow Not like contextual credit since we are always (regardless of context) interested in work of not hot f .

\rightarrow Heteroscedasticity, Time / History / Action dependent noise:
closer to RL problem! Information vs. Uncertainty

Nikolov, Kirschner, Berkenkamp ICML 2019

High-dims \Rightarrow Many vars \rightarrow structure \rightarrow subpace \rightarrow Murphy et al ICML