

Overview and Goals

Goal: AI systems that can learn and improve continually throughout their lives, just like humans and animals. Currently, deep learning (DL) requires a large amount of data which is costly and rigid, leading to a system that is unable to quickly adapt. We aim to fix this with a new learning paradigm based on Bayesian principles.

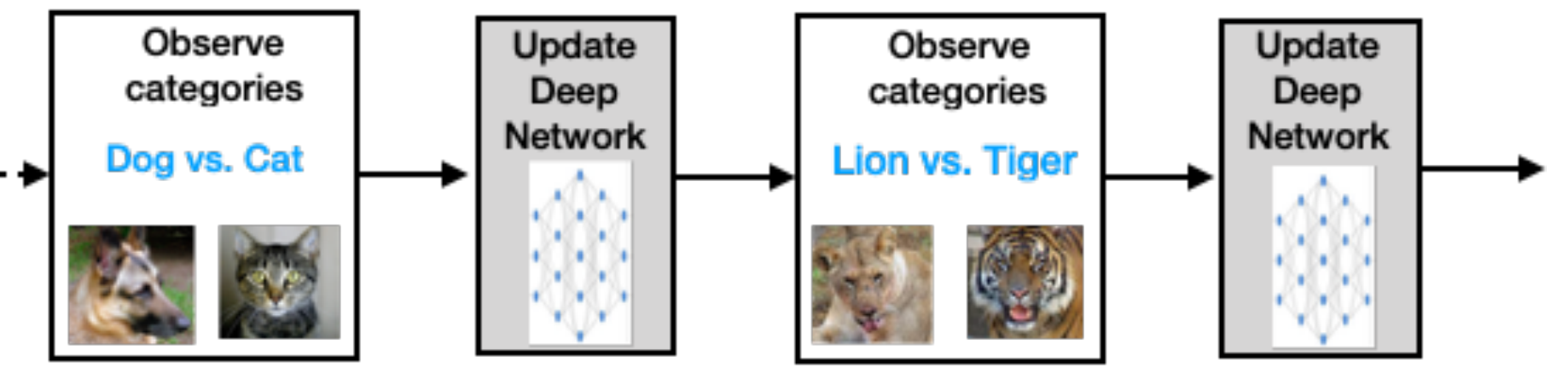
Summary of our research highlights in the year 2022 (paper number shown in green boxes)

- 1 We show a robust deep-learning method is related to Bayes and add uncertainty to it (top 75/5000 papers at ICLR2023)
- 2 We extend the Bayesian learning rule using Lie-groups, simplifying gradient computations and eliminating retractions.
- 3 We design a generally unimprovable procedure for relaxation time estimation in non-reversible Markov chains.
- 4 We simplify momentum-based Riemannian optimization over positive-semi-definite matrix submanifold.
- 5 We use low-rank matrix completion techniques to reconstruct partially-observed high-dimensional time series.
- 6-10 Analysis of MMD estimation, Neural processes, A Dataset for African Languages, Deviation inequalities, and more ...!

Standard Deep Learning



Deep Continual Learning

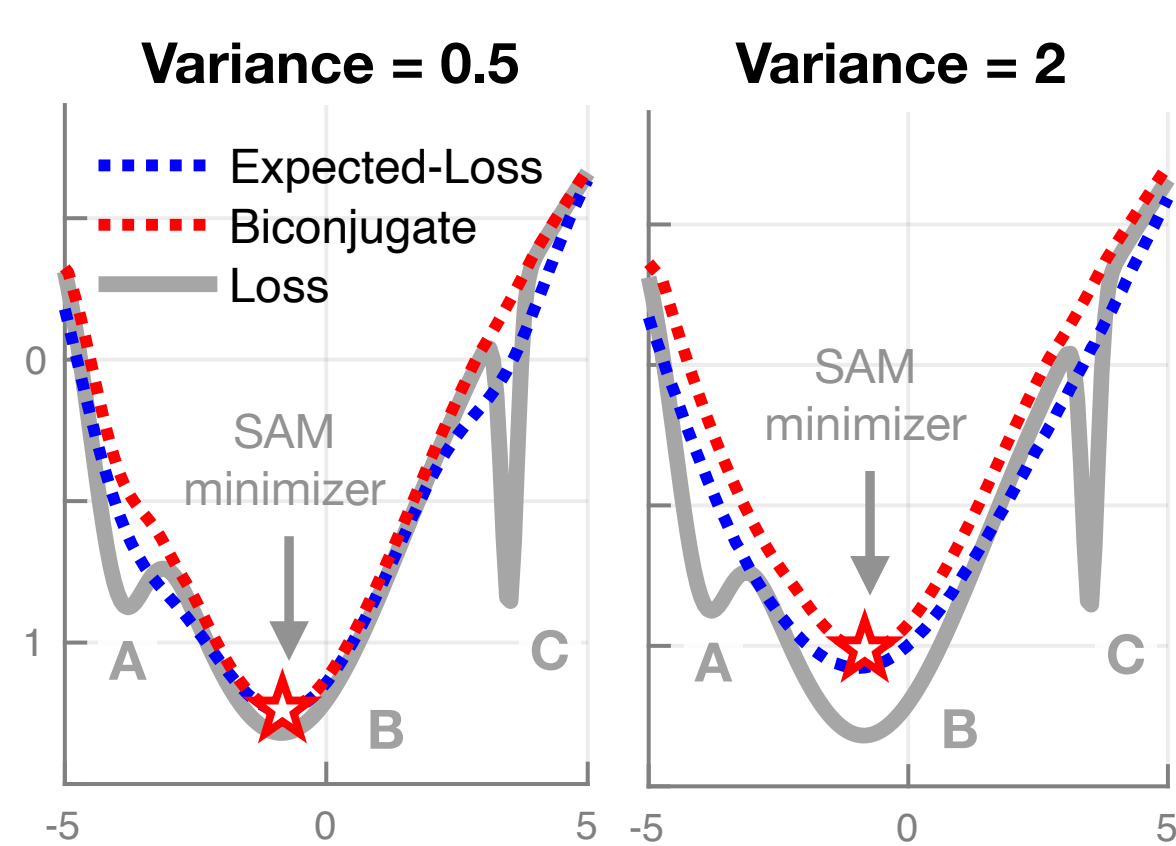
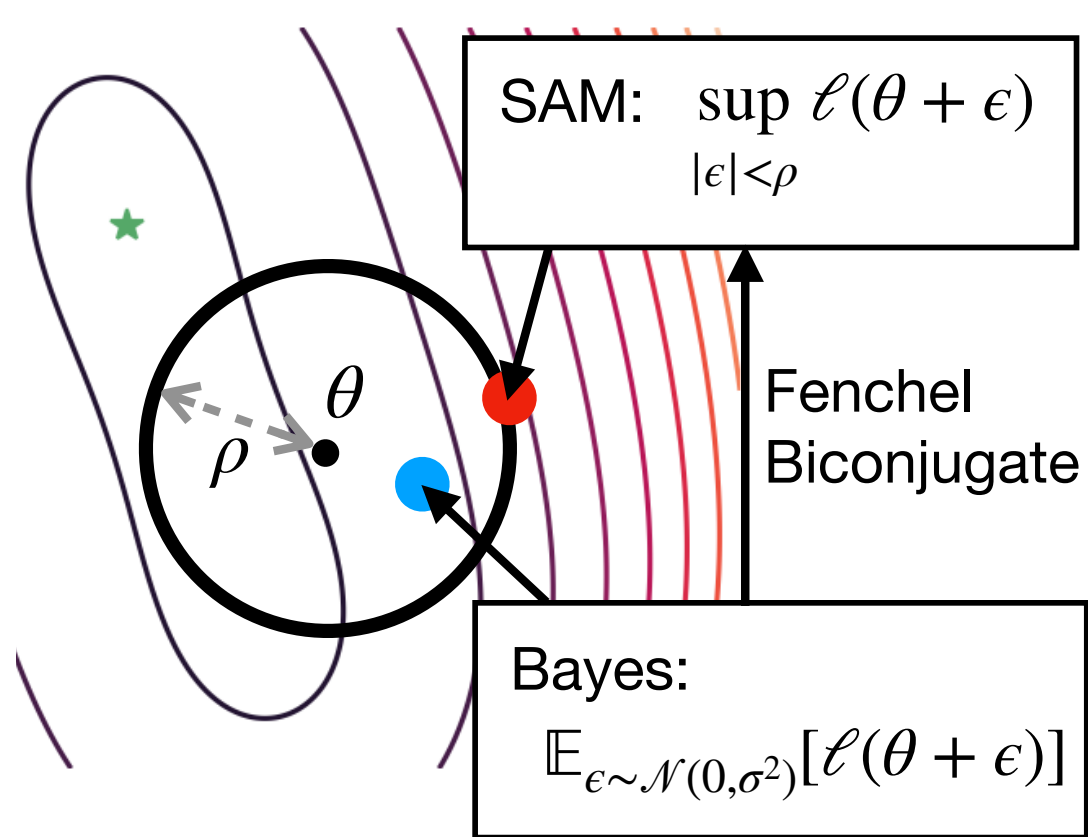


SAM as an Optimal Relaxation of Bayes

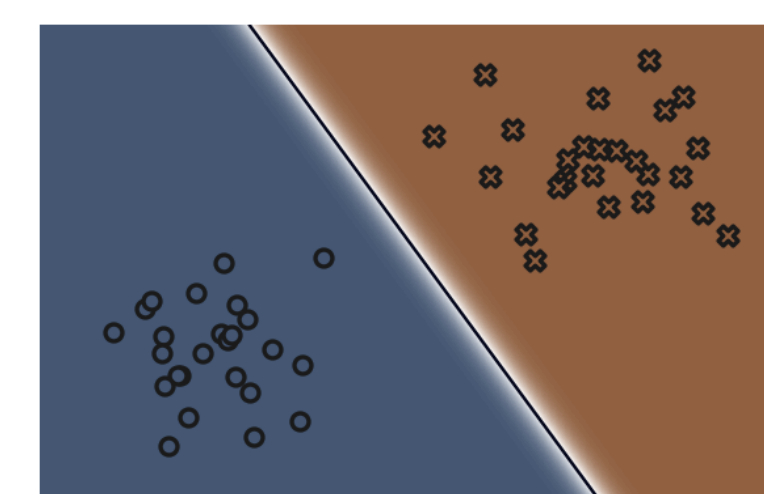
Problem: Sharpness-Aware Minimization (SAM) by Foret et al. improve significantly over SGD but the reasons behind its success are unclear.

Solution: We show that SAM is equivalent to an optimal relaxation of Bayes obtained by using Fenchel biconjugate (left figure). SAM can be seen as “smoothing” the objective using a “posterior variance” and always upper-bounding Bayes (right).

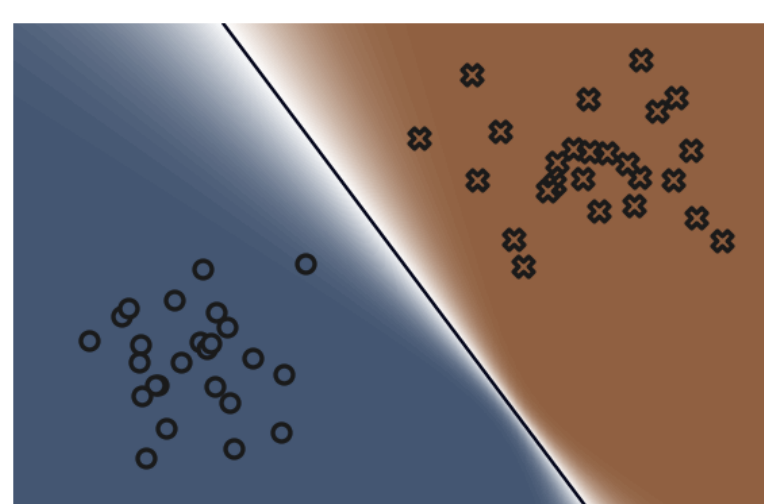
Our paper [1] is among top-5% of all accepted papers (75 out of 5000 submissions)



Contribution #2: Our Bayesian-SAM improves accuracy over SGD (by 8%), SAM-Adam (by 10%), and Adam (by 22%), but also improves AUROC (for CIFAR-100 using ResNet-20 with 270K params).



SAM is overconfident (narrow white region)



Bayesian-SAM fixes this (more blurry away from the data)

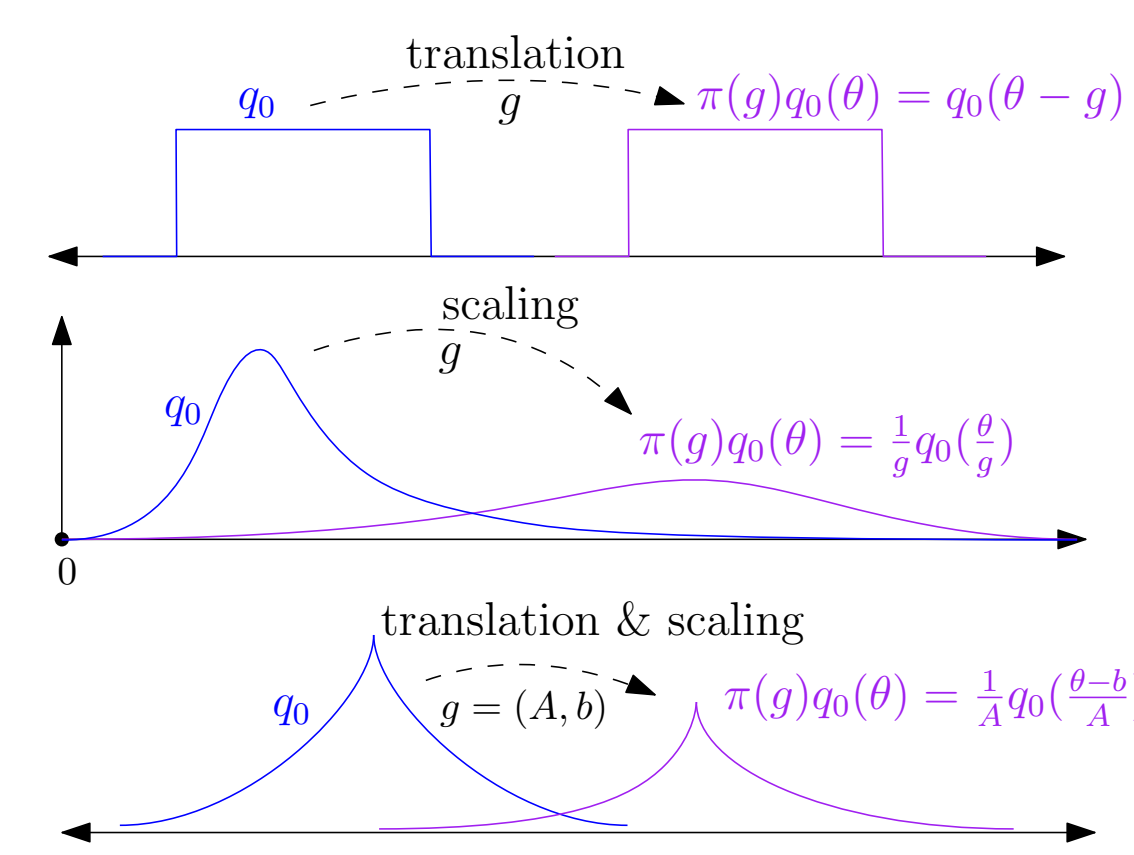
	Accuracy	AUROC
SGD	55.82(0.97) +8%	0.811(0.004)
SAM-SGD	58.58(0.59)	0.827(0.003)
SWAG	56.53(0.40)	0.814(0.004)
VOGN	59.83(0.75)	0.830(0.002)
Adam	39.73(0.97) +10%	0.775(0.004)
SAM-Adam	53.25(0.80) +22%	0.818(0.005)
bSAM (ours)	62.64(0.33)	0.841(0.004)

1. Moellenhoff, Khan, SAM as an Optimal Relaxation of Bayes, ICLR 2023 (oral, notable top-5%)

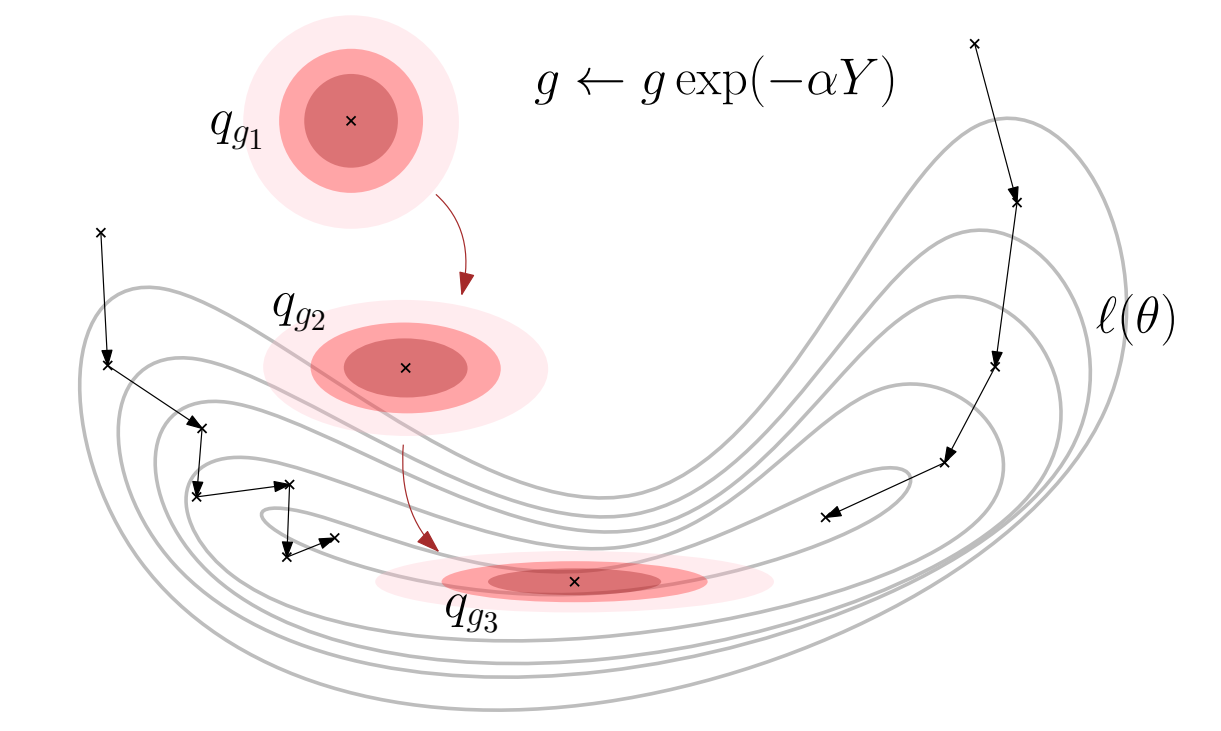
The Lie-Group Bayesian Learning Rule

Problem: Many popular algorithms can be derived from the Bayesian learning rule of Khan and Rue (2021) but the rule can be difficult to apply, e.g., gradients are difficult to compute, and steps can lead to invalid distributions (e.g., -ve variances).

Solution: We extend the rule by using Lie-groups which solves the above problems: gradients can always be obtained by reparameterizations, and steps always stay on the manifold. Fisher computation is also simplified and only need to be done once.



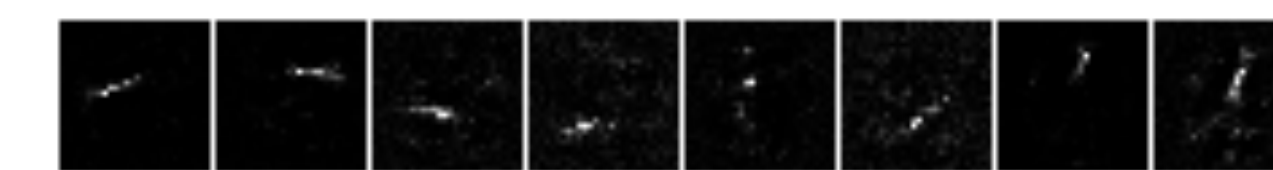
Conceptual algorithm:



Result: Reduces to the update [3] when applied to exp. fam. and linearized.

Additive $g \leftarrow g - \alpha \mathbb{E}_{q_g}[\nabla \theta \ell]$ **Multiplicative** $g_i \leftarrow g_i \exp(-\alpha (\mathbb{E}_{q_g}[\theta_i \partial_i \ell] - \tau))$ **Affine** $b_i \leftarrow b_i + \frac{c_x}{c_y} A_i \frac{\exp(-\alpha U) - 1}{U} V$, $A_i \leftarrow A_i \exp(-\alpha U)$

Multiplicative Learning (nll=0.058)



Additive Learning (nll=0.083)



Dataset	Metric	Affine (Alg. 3)			Additive (27, Alg. 1)		
		gaussian	laplace	uniform	gaussian	laplace	uniform
CIFAR-10	Acc. (↑)	91.53±0.10	91.87±0.04	91.60±0.05	91.28±0.11	91.14±0.12	91.07±0.08
	NLL (↓)	0.294±0.004	0.272±0.002	0.300±0.002	0.328±0.008	0.312±0.005	0.365±0.003
	ECE (↓)	0.036±0.001	0.029±0.001	0.040±0.001	0.045±0.001	0.039±0.001	0.052±0.001
CIFAR-100	Acc. (↑)	66.55±0.10	66.44±0.10	66.08±0.12	64.61±0.20	64.85±0.13	64.29±0.09
	NLL (↓)	1.255±0.005	1.247±0.006	1.288±0.007	1.390±0.008	1.359±0.006	1.437±0.004
	ECE (↓)	0.079±0.002	0.071±0.001	0.093±0.002	0.107±0.001	0.096±0.001	0.121±0.001
TinyImagNet	Acc. (↑)	51.13±0.16	51.36±0.14	51.19±0.12	49.62±0.15	49.73±0.18	49.34±0.14
	NLL (↓)	2.098±0.004	2.101±0.009	2.099±0.005	2.204±0.003	2.184±0.007	2.234±0.008
	ECE (↓)	0.070±0.002	0.065±0.001	0.076±0.001	0.099±0.002	0.089±0.001	0.107±0.001

2. Kiral, Moellenhoff, Khan, The Lie-Group Bayesian Learning Rule, AISTATS 2023.

Estimation of Relaxation Time in Non-Reversible Markov Chains

Problem: Estimate the relaxation time of an unknown, ergodic, non-reversible Markov chain from a single trajectory of length m started from an arbitrary state.

Application: This is useful in reinforcement learning, Markov Chain Monte Carlo diagnostics, and deriving generalization bounds with Markov-dependent data.

3. Wolfer and Kontorovich, Improved Estimation of Relaxation Time in Non-reversible Markov Chains, arXiv 2022.

Contributions:

- Prove there exists an estimator \hat{t}_{rel} for t_{rel} such that for any ergodic chain, when $m = \tilde{O}(t_{rel}/\pi_*)$ where π_* is the minimum stationary probability, it holds w.h.p. that $|\hat{t}_{rel}/t_{rel} - 1| < 5$. Upper bound matches known lower bounds, thus **generally unimprovable**.
- Design algorithm with **fully empirical confidence intervals** that decay in $\tilde{O}(1/\sqrt{m})$.

Other Works

Simplify momentum-based Riemannian optimization: We consider the specific case of a submanifold containing symmetric positive-definite matrices. The method uses a generalized version of local coordinates which “trivializes” the Fisher matrix.

4. Lin, Duruisseaux, Leok, Nielsen, Khan, Schmidt, Practical Structured Riemannian Optimization with Momentum by using Generalized Normal Coordinates, NeurIPS 2022 Workshop on Symmetry and Geometry in Neural Representation

High-dimensional time series completion: We use low-rank matrix completion techniques to reconstruct partially observed high-dimensional time series and show that periodicity or smoothness can even lead to faster rates than in the independent setting.

5. Alquier, Marie, Rosier, Tight risk bound for high dimensional time series completion, EJS 2022

Finite sample properties of parametric MMD estimation: We tackle the problem of universal estimation using a minimum distance estimator based on Maximum Mean Discrepancy, and we show its robustness to both dependence and presence of outliers.

6. Chérif-Abdellatif, Alquier, Finite Sample Properties of Parametric MMD Estimation: Robustness to Misspecification and Dependence, Bernoulli 2022

Improving Neural Process: We improve test-time inference for Neural Processes by incorporating and exploiting graphical-model structure among context points.

7. Tailor, Khan, Nalisnick, Exploiting Inference Structure in Neural Processes, UAI 2022 Workshop on Tractable Probabilistic Modeling

Name-Entity Recognition (NER) dataset for Sub-Saharan African languages: We create the largest human-annotated dataset called MasakhaNER 2.0, and analyze features that contribute to cross-lingual transfer, giving large gains for 0-shot learning.

8. Buzaaba with many others, MasakhaNER 2.0: Africa-centric Transfer Learning for Named Entity Recognition, EMNLP 2022

Deviation inequalities for stochastic approximation by averaging: We establish deviation inequalities for separately Lipschitz functions of Markov chains belonging to a certain class we define, and which includes models of stochastic approximation by averaging and non-averaging.

9. Fan, Alquier, Doukhan, Deviation inequalities for stochastic approximation by averaging, SPA 2022

Prioritization of minibatches: We give empirical support for the hypothesis that improving calibration can help in prioritizing minibatches during training.

10. Tata, Gudur, Chennupati, Khan, Can calibration improve sample prioritization, NeurIPS 2022 Workshop on Has It Trained Yet