



# SAM as an optimal approximation of Bayes

#### Mohammad Emtiyaz Khan RIKEN Center for AI Project, Tokyo http://emtiyaz.github.io



Summary of research at <a href="https://emtiyaz.github.io/papers/symposium\_2021.pdf">https://emtiyaz.github.io/papers/symposium\_2021.pdf</a> Slides available at <a href="https://emtiyaz.github.io/papers/Nov24\_2022\_Ubath.pdf">https://emtiyaz.github.io/papers/Nov24\_2022\_Ubath.pdf</a>

# Al that learn like humans

Quickly adapt to learn new skills, throughout their lives

# Human Learning at the age of 6 months.



# Converged at the age of 12 months



Transfer skills at the age of 14 months



# Fail because too slow to adapt



https://www.youtube.com/watch?v=TxobtWAFh8o The video is from 2017

# Al that learn like humans

Quickly adapt to learn new skills, throughout their lives

#### Thomas Moellenhoff

# **Robust Deep-Learning**

- Sharpness-Aware Minimization (SAM)[1]
  - Huge improvements over SGD/Adam
  - Now used to train all sorts of models
  - Improves test accuracy for trained neural networks (ResNets [1,2,3], Vision-Transformers [4], Language Models [5], ...)
  - Also improves robustness [3], calibration [2], interpretability [4], transfer-learning [4], compressibility [5], federated learning [6]
- Why does it work, and how to improve it?
- SAM as an "optimal" relaxation of Bayes [7]
- 1. Foret et al. Sharpness-aware minimization for efficiently improving generalization. ICLR 2021.
- 2. Wu et al. Adversarial Weight Perturbation Helps Robust Generalization . *NeurIPS 2020.*
- 3. Zheng et al. Regularizing Neural Networks via Adversarial Model Perturbation . CVPR 2021.
- 4. Chen et al. When Vision Transformers Outperform ResNets without Pretraining or Strong Data Augmentations. *ICLR 2022.*
- 5. Na et al. Train Flat, Then Compress: SAM Learns More Compressible Models. EMNLP 2022.
- 6. Qu et al. Generalized Federated Learning via Sharpness Aware Minimization. ICML 2022.
- 7. Moellenhoff and Khan, SAM as an optimal relaxation of Bayes, arXiv, 2022.

# Flat Minima in DL

- Empirical [1] and theoretical [2] evidence suggest that finding "flat minima" is desirable
  - Several ways to bias the learning towards flatter minima, e.g., in SGD via batch-size or learning rate



- 1. Jiang et al., Fantastic generalization measures and where to find them, ICLR 2020.
- 2. Dzuigate and Roy, Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data, UAI 2017.

### Flat Minima via Learning Rate in SGD





#### Flat Minima via Adversarial Weight-Perturbation



## Flat minima via Bayesian Learning



VS

min  $\ell(\theta)$ 

Gaussian approximation

 $\min_{q \in \mathcal{Q}} \mathbb{E}_{q(\theta)}[\ell(\theta)] - \mathcal{H}(q)$ Entropy

 First term"smooths" the loss by weight perturbation

 $\mathbb{E}_{\mathcal{N}(\epsilon|0,\sigma^2)}[\ell(m+\epsilon)]$ 

- Similar mechanism to SAM by it is the "expected-loss" instead of "max-loss"
- An advantage of Bayes is that the second term can be used to estimate  $\sigma$
- Foret et al. Sharpness-Aware Minimization for Efficiently Improving Generalization, ICLR, 2021
   Smith et al., On the Origin of Implicit Regularization in Stochastic Gradient Descent, ICLR, 2021

#### SAM as an Optimal relaxation of Bayes



# SAM as a relaxation of Bayes

SAM (red star) upper bounds the Bayesian  $\mathbb{E}_q[\ell]$ 



### **Optimal relaxation: Fenchel Conjugate**

SAM minimizes the best-concave upper bound to  $\mathbb{E}_q[\mathcal{E}]$  wrt the mean, while keeping variance fixed.



# **Legendre-Fenchel Transform**

Fenchel conjugate

$$f^*(v) = \sup_{u \in U} \langle v, u \rangle - f(u)$$

Fenchel biconjugate



Our bound is obtained in the 2D space of "expectation parameters" of Gaussians win mean  $\omega$  and variance v Expectation parameter  $\mu = (\omega, \omega^2 + v)$ Variance Mean  $f(\mu) = \mathbb{E}_{\theta \sim q_{\mu}}[-\ell(\theta)]$ 0 -Ő 20  $\omega^2 + v$ 

The bound in the previous slide is in  $(\omega, v)$  parameterization

## **Relaxed Bayesian Objective**

• The relaxed objective

 $\mathcal{L}(q_{\mu}) \geq \mathcal{L}_{\text{relaxed}}(q_{\mu}) = f^{**}(\mu) - \mathbb{D}_{\text{KL}}[q_{\mu}(\theta) || p(\theta)].$ 

- This is a difference of convex objective, for which there is a **dual problem** [1] Natural parameter  $\max_{\mu} f^{**}(\mu) - g(\mu) = \max_{\lambda} \left[ -f^{*}(\lambda) + g^{*}(\lambda) \right]$
- In our case, the dual problem can be rewritten in the following form:

$$\mathcal{E}_{\text{relaxed}}(\mathbf{m},\sigma;\delta') = \left[\sup_{\boldsymbol{\epsilon}} \ell(\mathbf{m}+\boldsymbol{\epsilon}) - \frac{1}{2\sigma^2} \|\boldsymbol{\epsilon}\|^2\right] + \underbrace{\frac{\delta'}{2} \|\mathbf{m}\|^2 - \frac{P}{2} \log(\sigma^2 \delta')}_{= -\log \mathcal{Z}},$$

### **Our Result: Relaxed Bayes is SAM!**

$$\begin{split} \mathcal{E}_{\text{relaxed}}(\mathbf{m},\sigma;\delta') &= \left[ \sup_{\boldsymbol{\epsilon}} \,\ell(\mathbf{m}+\boldsymbol{\epsilon}) - \frac{1}{2\sigma^2} \|\boldsymbol{\epsilon}\|^2 \right] + \underbrace{\frac{\delta'}{2} \|\mathbf{m}\|^2 - \frac{P}{2} \log(\sigma^2 \delta')}_{= -\log \mathcal{Z}}, \\ &= -\log \mathcal{Z} \end{split}$$

**Theorem 2.** For every  $(\rho, \delta)$ , there exist  $(\sigma, \delta')$  such that

 $\arg\min_{\boldsymbol{\theta}\in\Theta} \ \mathcal{E}_{SAM}(\boldsymbol{\theta};\rho,\delta) = \arg\min_{\mathbf{m}\in\Theta} \ \mathcal{E}_{relaxed}(\mathbf{m},\sigma;\delta').$ 

Essentially, for fixed variance, we can always recover SAM's solution by minimizing the relaxedbases objective with respect to the mean

# **Improving SAM**

We can obtain the variance by optimizing the relaxed Bayes objective. Below is the Fenchel biconjugate where a covariance  $\Sigma$  is used

$$-f^{**}(\boldsymbol{\mu}) = \min_{\mathbf{m}\in\mathbb{R}^{P},\mathbf{b}\in\mathbb{R}^{P}_{+}} \left[ \sup_{\boldsymbol{\epsilon}\in\mathbb{R}^{P}} \ell(\mathbf{m}+\boldsymbol{\epsilon}) - \frac{1}{2} \|\boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{\epsilon}\|^{2} \right] + \mathbb{E}_{\boldsymbol{\theta}\sim q_{\boldsymbol{\mu}}} \left[ \frac{1}{2} \|\boldsymbol{\Sigma}^{-\frac{1}{2}}(\boldsymbol{\theta}-\mathbf{m})\|^{2} \right]$$
  
Covariance

We use the Bayesian learning rule (BLR) [1] to do that, because it gives an Adam-Style objective

$$\begin{array}{c} \boldsymbol{\lambda} \leftarrow (1-\alpha)\boldsymbol{\lambda} + \alpha \left[ \nabla f^{**}(\boldsymbol{\mu}) + \boldsymbol{\lambda}_0 \right] \\ \uparrow & \uparrow \\ \text{Natural parameter} & \text{Natural gradient} & \text{Prior} \end{array}$$

# **Bayesian-SAM**

An Adam-style algorithm, derived using the BLR, where "perturbation-size" is adjusted by using  $\sigma^2$  (or *s*)

SAM with RMSprop

 $g_{1} \leftarrow \hat{\nabla}\ell(\theta)$   $\epsilon \leftarrow \rho \frac{g_{1}}{\|g_{1}\|}$   $g \leftarrow \hat{\nabla}\ell(\theta + \epsilon)$   $s \leftarrow (1 - \rho)s + \rho g^{2}$  $\theta \leftarrow \theta - \alpha(\sqrt{s} + \delta)^{-1}g$  SAM with BLR

$$g_{1} \leftarrow \hat{\nabla}\ell(\theta), \text{ where } \theta \sim \mathcal{N}(m, \sigma^{2})$$
  

$$\epsilon \leftarrow \rho \frac{g_{1}}{s}$$
  

$$g \leftarrow \hat{\nabla}\ell(\theta + \epsilon)$$
  

$$s \leftarrow (1 - \rho)s + \rho \sqrt{s}|g_{1}|$$
  

$$m \leftarrow m - \alpha(s + \delta)^{-1} \nabla_{\theta}\ell(\theta)$$
  

$$\sigma^{2} \leftarrow (s + \delta)^{-1}$$

Foret et al. Sharpness-Aware Minimization for Efficiently Improving Generalization, ICLR, 2021
 Moellenhoff and Khan, SAM as an optimal relaxation of Bayes, https://arxiv.org/abs/2210.01620, 2022

Bayes



#### Improving "overconfident" SAM



Bayesian SAM



SAM



# Results on Resnet-20 (200K params) and ResNet-18 (11M params)

Model / Dataset	Method	Accuracy † (higher is better)	NLL↓ (lower is better)	ECE $\downarrow$ (lower is better)	AUROC ↑ (higher is better)
ResNet-20-FRN / SAM-SGD CIFAR-10 Adam SAM-Adam bSAM (ours)		$\begin{array}{c} 91.68_{(0.26)}\\ \textbf{92.29}_{(0.39)}\\ 89.97_{(0.27)}\\ 91.57_{(0.21)}\\ \textbf{92.16}_{(0.16)}\end{array}$	$\begin{array}{r} 0.29_{(0.008)}\\ 0.25_{(0.004)}\\ 0.41_{(0.021)}\\ 0.26_{(0.004)}\\ \hline 0.23_{(0.003)}\end{array}$	$\begin{array}{c} 0.0397_{(0\ 002)}\\ 0.0266_{(0\ 003)}\\ 0.0610_{(0\ 002)}\\ 0.0329_{(0\ 002)}\\ 0.0057_{(0\ 002)}\end{array}$	$\begin{array}{c} 0.915_{(0.002)} \\ 0.920_{(0.003)} \\ 0.900_{(0.003)} \\ 0.918_{(0.001)} \\ 0.925_{(0.001)} \end{array}$
SGD ResNet-20-FRN / SAM-SGD CIFAR-100 Adam SAM-Adam bSAM (ours)		$\begin{array}{r} 66.48_{(0.10)} \\ 67.27_{(0.22)} \\ 61.76_{(0.67)} \\ 65.34_{(0.32)} \\ 68.22_{(0.14)} \end{array}$	$1.20_{(0.007)}$ $1.19_{(0.011)}$ $1.66_{(0.049)}$ $1.23_{(0.012)}$ $1.10_{(0.012)}$	$\begin{array}{c} 0.0524_{(0.004)}\\ 0.0481_{(0.001)}\\ 0.1582_{(0.006)}\\ 0.0166_{(0.003)}\\ 0.0258_{(0.003)}\end{array}$	$\begin{array}{c} 0.846_{(0.002)} \\ 0.848_{(0.002)} \\ 0.826_{(0.003)} \\ 0.847_{(0.002)} \end{array}$
ResNet-20-FRN / SAM-SGD TinyImageNet Adam SAM-Adam bSAM (ours)		$\begin{array}{c} 52.01_{(0.36)}\\ 52.25_{(0.26)}\\ 49.04_{(0.38)}\\ 51.17_{(0.45)}\\ 52.90_{(0.35)}\end{array}$	$\begin{array}{r} 1.98_{(0.007)}\\ 1.97_{(0.013)}\\ 2.14_{(0.024)}\\ 2.02_{(0.014)}\\ 1.94_{(0.009)}\end{array}$	$\begin{array}{c} 0.0330_{(0.032)}\\ \textbf{0.0155}_{(0.032)}\\ 0.0502_{(0.034)}\\ 0.0460_{(0.034)}\\ \textbf{0.0199}_{(0.033)}\end{array}$	0.832 <sub>(0.002)</sub> 0.827 <sub>(0.005)</sub> 0.820 <sub>(0.004)</sub> 0.828 <sub>(0.004)</sub> 0.831 <sub>(0.001)</sub>
ResNet-18 / CIFAR-10	SGD SAM-SGD b <b>SAM (ours)</b>	$94.76_{(0.11)} \\95.72_{(0.14)} \\96.15_{(0.08)}$	$\begin{array}{c} 0.21_{(0.006)}\\ 0.14_{(0.004)}\\ 0.12_{(0.002)}\end{array}$	$\begin{array}{c} 0.0304_{(0.001)}\\ 0.0134_{(0.001)}\\ 0.0049_{(0.001)}\end{array}$	$\begin{array}{c} 0.926_{(0.006)}\\ 0.949_{(0.003)}\\ 0.954_{(0.001)}\end{array}$
ResNet-18 / CIFAR-100	SGD SAM-SGD bSAM (ours)	$76.54_{(0.26)}$ $78.74_{(0.19)}$ $80.22_{(0.28)}$	0.98 <sub>(0.007)</sub> 0.79 <sub>(0.007)</sub> <b>0.70</b> <sub>(0.008)</sub>	$\begin{array}{c} 0.0501_{(0.002)}\\ 0.0445_{(0.002)}\\ \textbf{0.0311}_{(0.003)}\end{array}$	$\begin{array}{c} 0.869_{(0.003)} \\ 0.887_{(0.003)} \\ \textbf{0.892}_{(0.003)} \end{array}$

### Bayesian SAM is less sensitive to hyper-parameters



#### **The Bayes-Duality Project**

Toward AI that learns adaptively, robustly, and continuously, like humans



Emtiyaz Khan

Research director (Japan side)

Approx-Bayes team at RIKEN-AIP and OIST Julyan Arbel

Research director (France side)

Statify-team, Inria Grenoble Rhône-Alpes

#### Kenichi Bannai

Co-PI (Japan side)

Math-Science Team at RIKEN-AIP and Keio University Rio Yokota

(Japan side)

Tokyo Institute of Technology

Received total funding of around USD 3 million through JST's CREST-ANR and Kakenhi Grants.

## **Approximate Bayesian Inference Team**

#### https://team-approx-bayes.github.io/



Emtiyaz Khan Team Leader



Pierre Alguier Research Scientist



<u>Hugo Monzón</u> Maldonado Postdoc



Happy Buzaaba Postdoc



Erik Daxberger Remote Collaborator University of Cambridge



Paul Chang Remote Collaborator Aalto University



Gian Maria Marconi Postdoc



Thomas Möllenhoff Postdoc



Lu Xu Postdoc



Jooyeon Kim Postdoc



<u>Alexandre Piché</u> Remote Collaborator *MIL*A



Ang Mingliang Remote Collaborator National University of Singapore



Geoffrey Wolfer Postdoc



Wu Lin PhD Student University of British Columbia



Peter Nickl Research Assistant



Dharmesh Tailor Remote Collaborator University of Amsterdam