

# The Bayesian Learning Rule

Mohammad **Emtiyaz** Khan

RIKEN Center for AI Project, Tokyo

<http://emtiyaz.github.io>



# **AI 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 quick to adapt

## TayTweets: Microsoft AI bot manipulated into being extreme racist upon release

Posted Fri 25 Mar 2016 at 4:38am, updated Fri 25 Mar 2016 at 9:17am



TayTweets is programmed to converse like a teenage girl who has "zero chill", according to Microsoft. (Twitter/TayTweets)

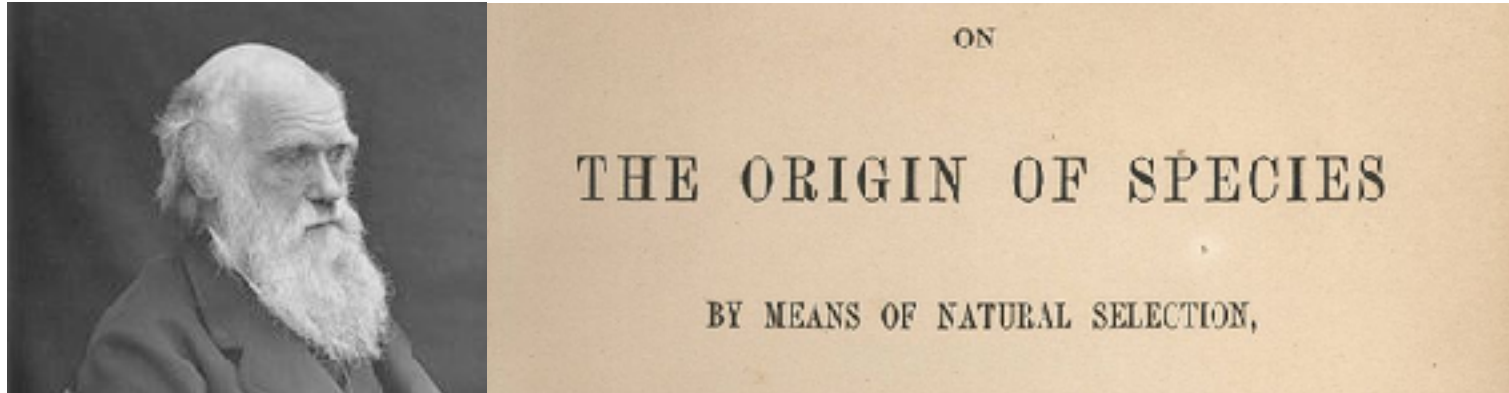
# Failure of AI in “dynamic” setting

Robots need quick adaptation to be deployed  
(for example, at homes for elderly care)



# **AI that learn like humans**

Quickly adapt to learn new skills, throughout their lives



# The Origin of Algorithms

What are the common principles behind popular algorithms?

# Principles of “good” algorithms?

- Why Bayes?
- Information Geometry of Bayes
  - To unify/generalize/improve learning-algorithms
  - Optimize for “posterior approximations”
- Bayesian Learning rule (BLR)
  - Derive many algorithms from optimization, deep learning, and Bayesian inference
- Natural Gradients are Everywhere!

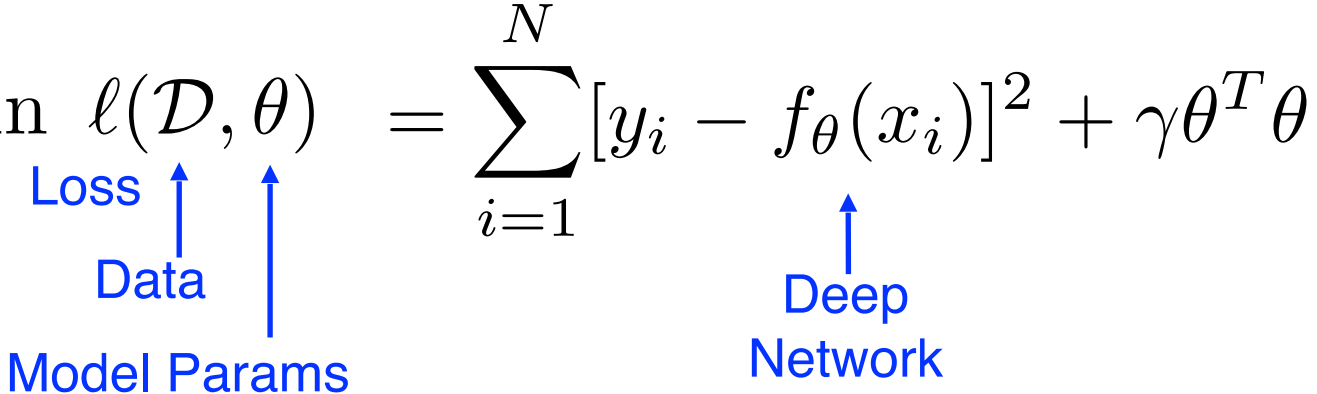


# Why Bayes?

Nasty data, adaptation, uncertainty  
estimation, reducing overfitting,  
model selection

# Principle of Trial-and-Error

Frequentist: Empirical Risk Minimization (ERM) or Maximum Likelihood Principle, etc.

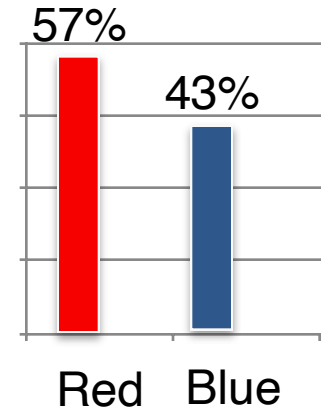
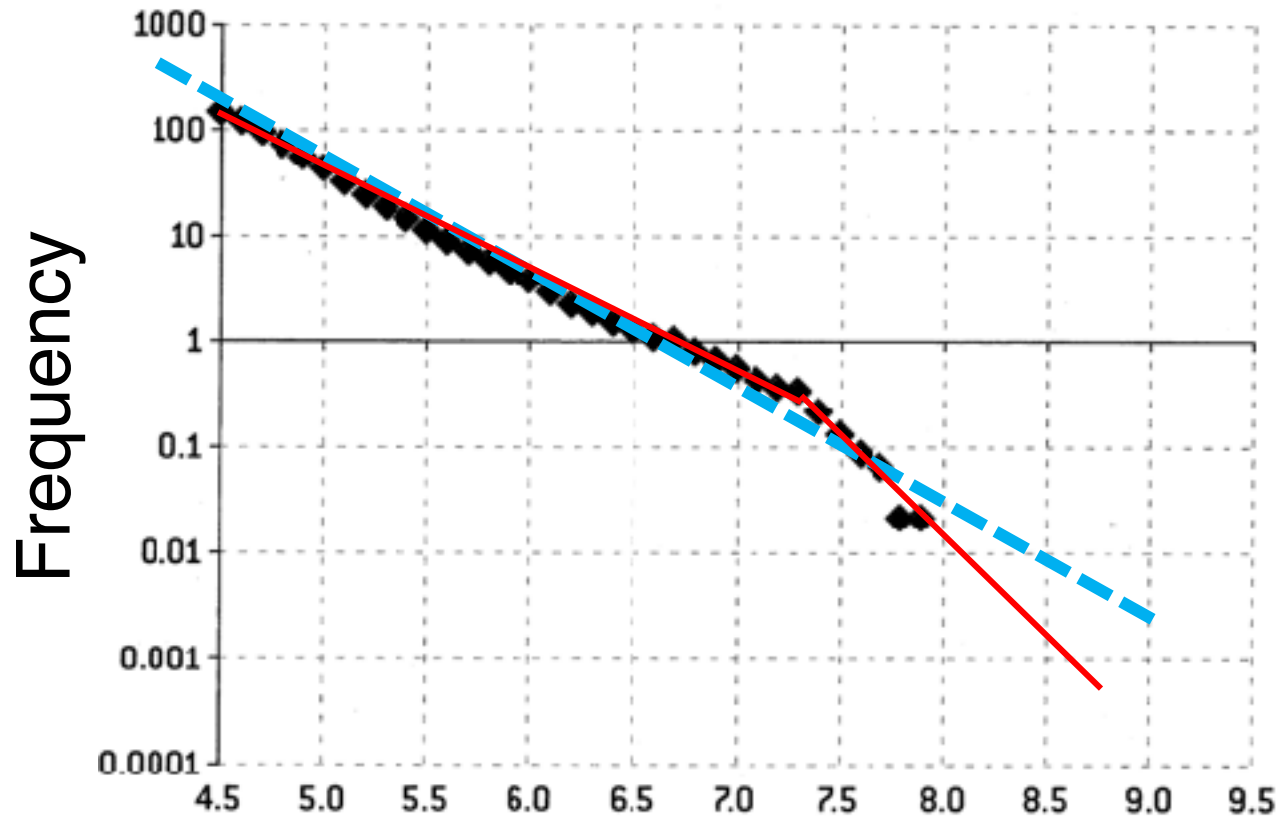
$$\min_{\theta} \ell(\mathcal{D}, \theta) = \sum_{i=1}^N [y_i - f_{\theta}(x_i)]^2 + \gamma \theta^T \theta$$


Loss ↑  
Data ↑  
Model Params ↑  
Deep Network ↑

Deep Learning Algorithms:  $\theta \leftarrow \theta - \rho H_{\theta}^{-1} \nabla_{\theta} \ell(\theta)$

Scales well to large data and complex model, and very good performance in practice.

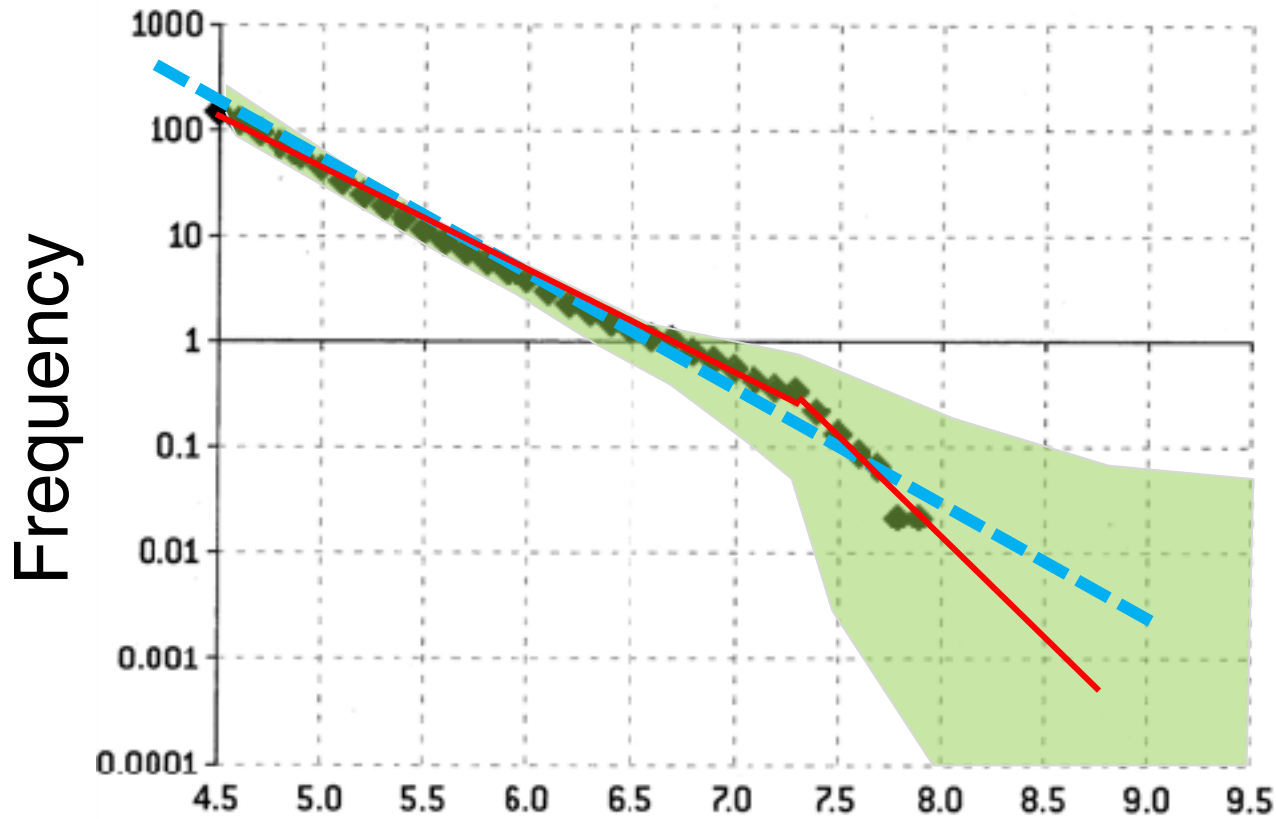
# Example: Which is a Better Fit?



More data  $\longrightarrow$  Less data  
Magnitude of Earthquake

Red is more  
risky than  
the blue

# Example: Which is a Better Fit?



More data  $\longrightarrow$  Less data  
Magnitude of Earthquake

Uncertainty:  
“What the  
model does  
not know”

Choose less  
risky options!

Avoid data  
bias with  
uncertainty!

# Bayesian Principles

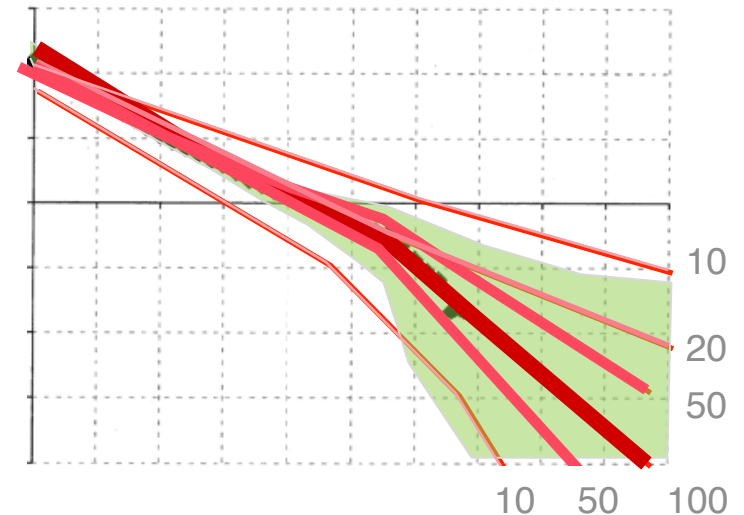
1. Sample  $\theta \sim p(\theta)$  prior

2. Score  $p(\mathcal{D}|\theta) = \prod_{i=1}^N p(y_i | f_{\theta}(x_i))$  Likelihood

3. Normalize

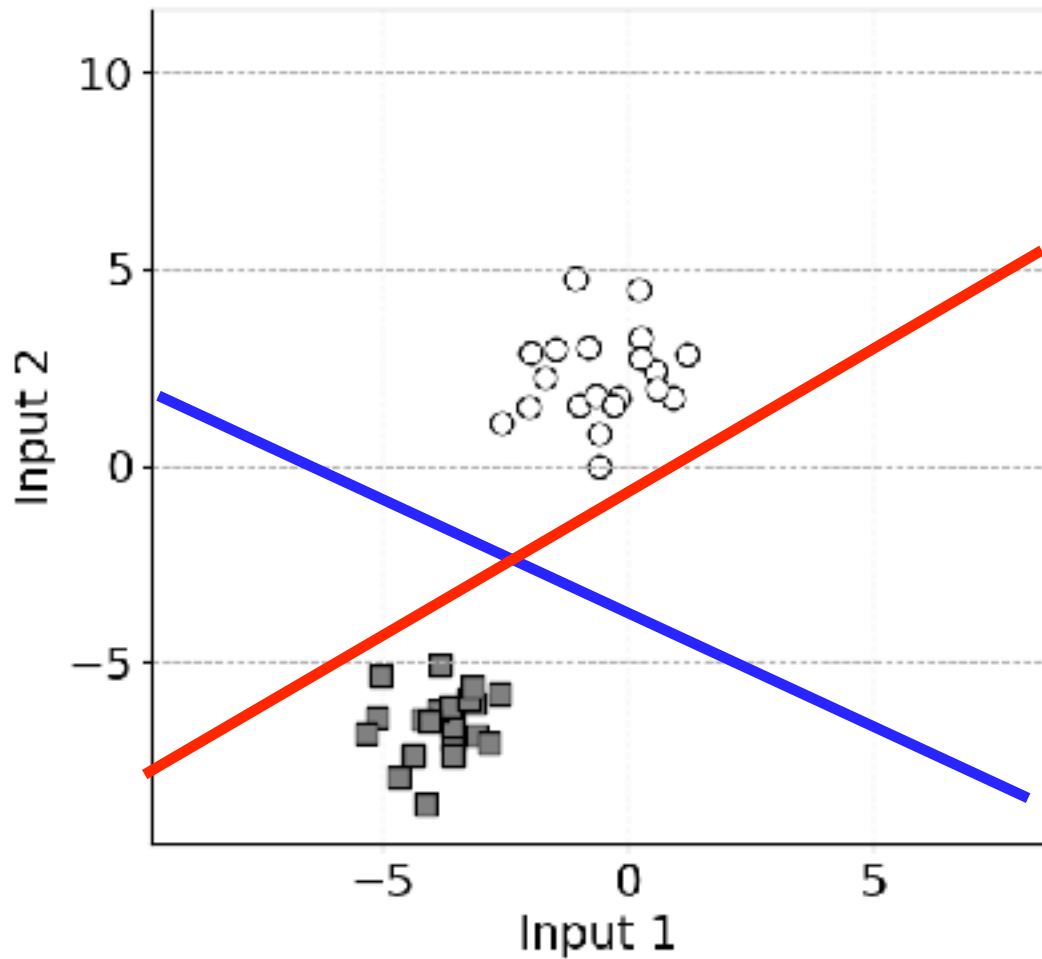
Posterior Likelihood x Prior

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{\int p(\mathcal{D}|\theta)p(\theta)d\theta}$$



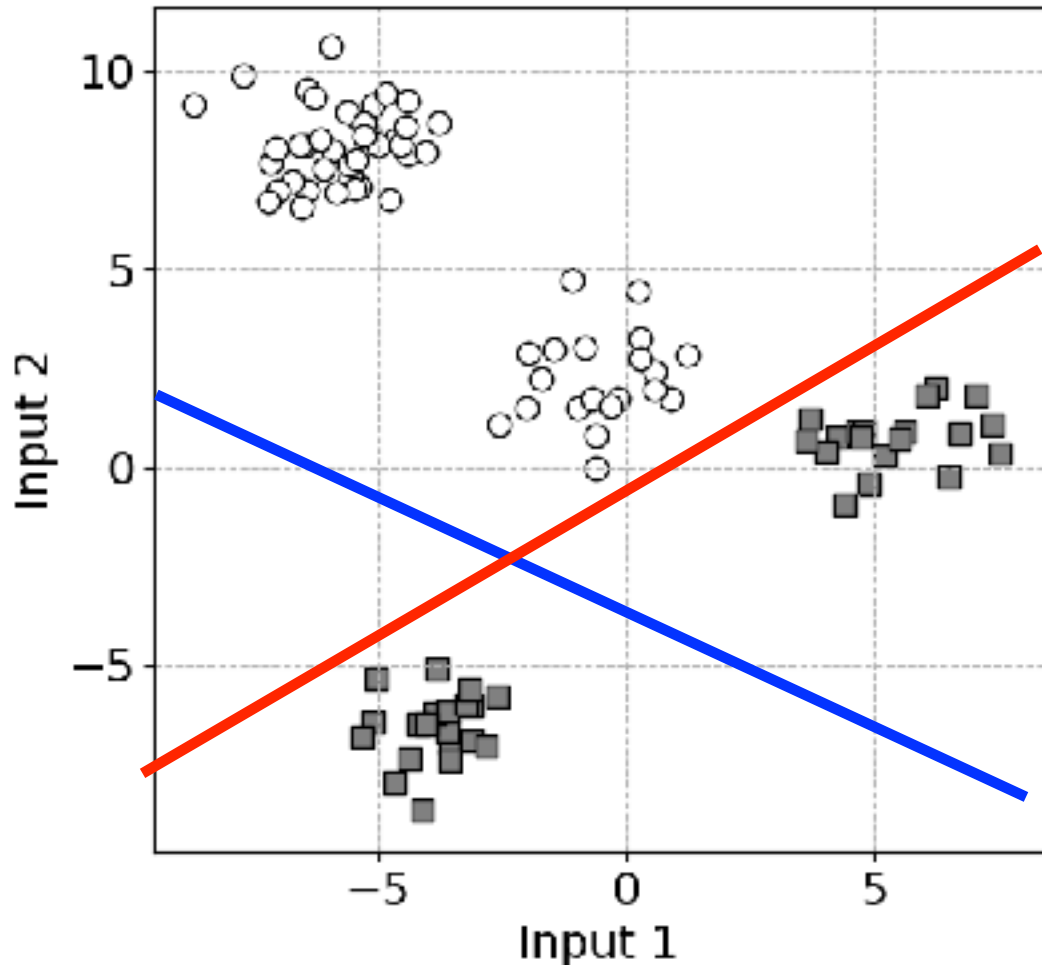
A global method: Integrates over all models  
Does not scale to large problem

# Which is a good classifier?





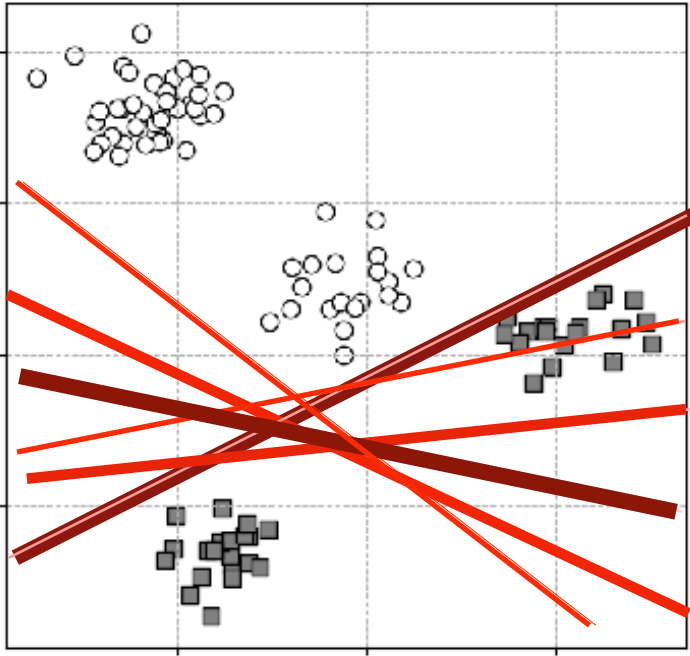
# Which is a good classifier?



Misclassified by the red line, but not by the blue

What you don't know now, can hurt you later  
**“Uncertainty matters”**

# Bayesian Principles



(1) Keep your options open

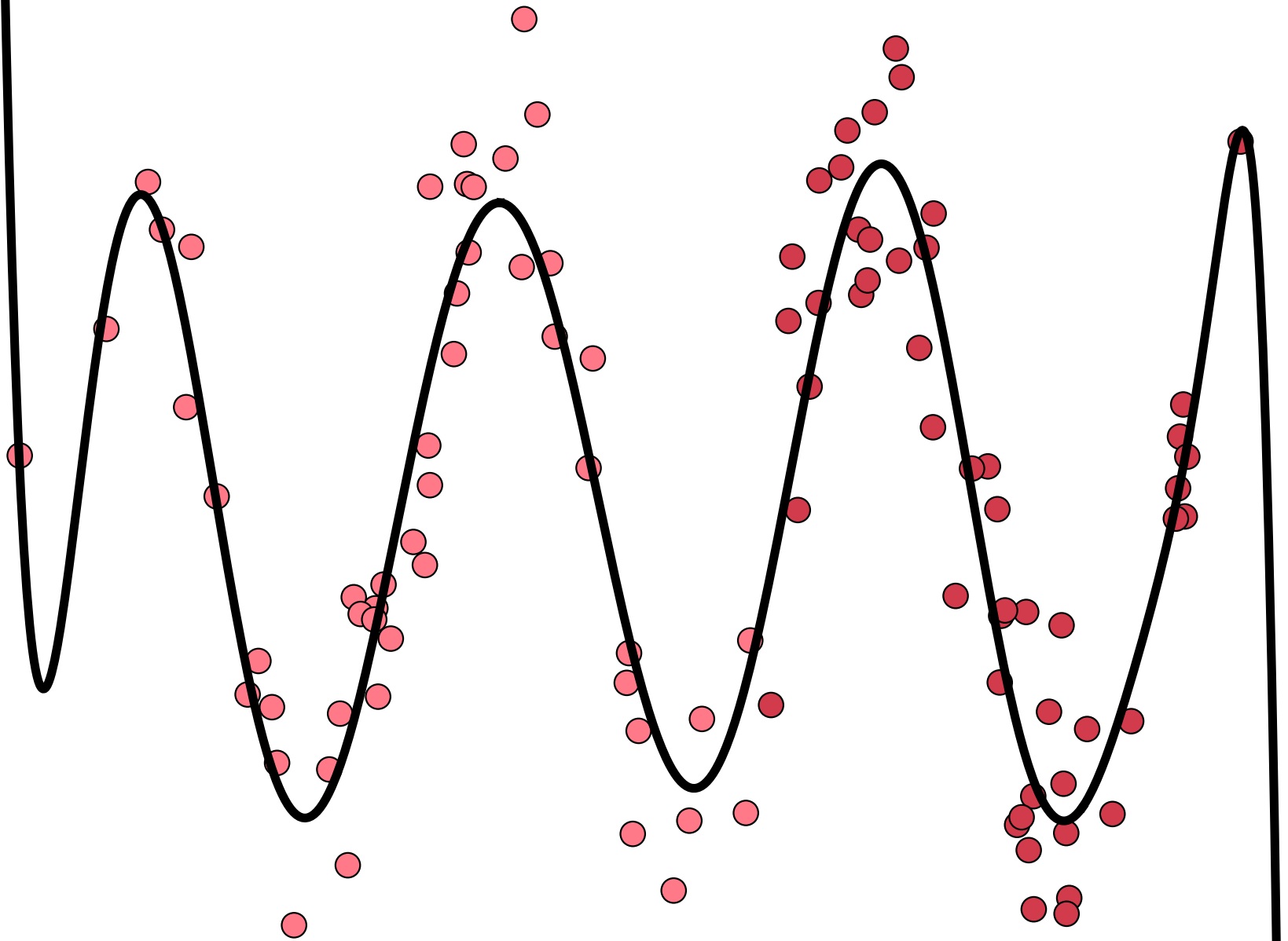
$$p(\theta|\mathcal{D}_1) = \frac{p(\mathcal{D}_1|\theta)p(\theta)}{\int p(\mathcal{D}_1|\theta)p(\theta)d\theta}$$

(2) Revise with new evidence

$$p(\theta|\mathcal{D}_2, \mathcal{D}_1) = \frac{p(\mathcal{D}_2|\theta)p(\theta|\mathcal{D}_1)}{\int p(\mathcal{D}_2|\theta)p(\theta|\mathcal{D}_1)d\theta}$$

Similar ideas in sequential/online decision-making (uncertainty/randomization). **Computation is infeasible.**

# Bayesian Linear Regression (polynomials of degree 15)



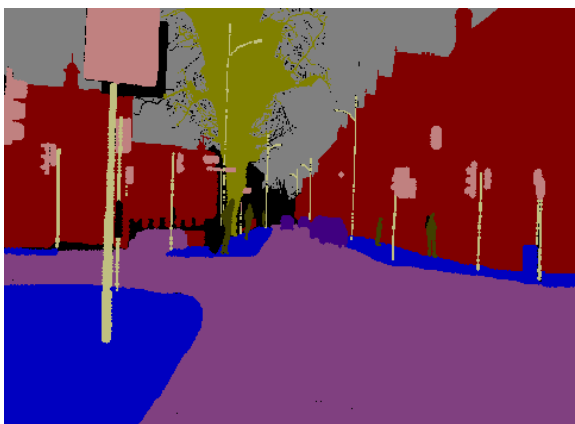
(By Roman Bachmann)

# Uncertainty Estimates for Image Segmentation

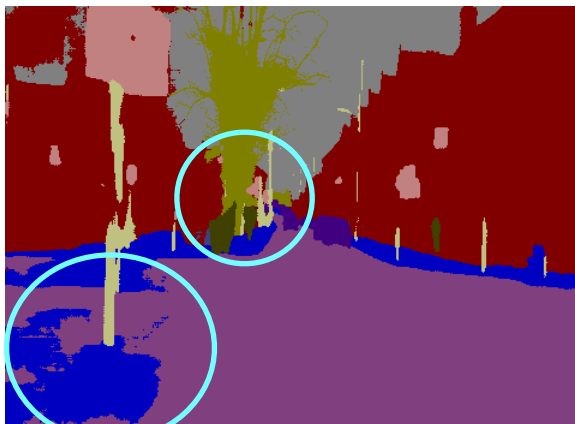
Image



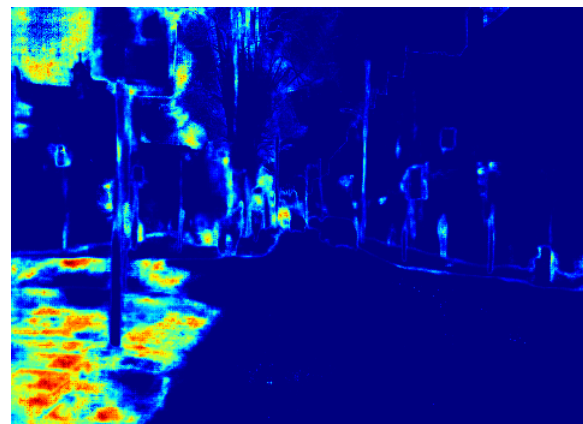
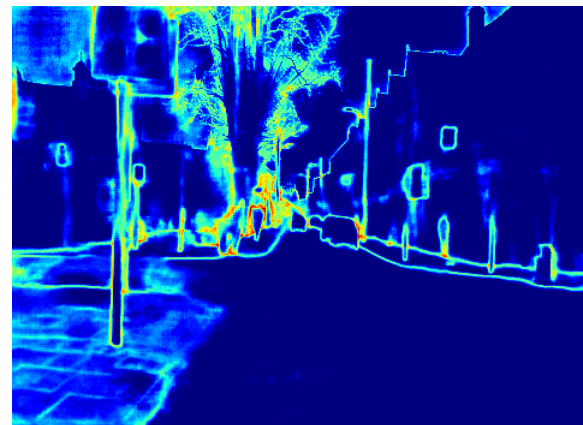
True Segments



Prediction



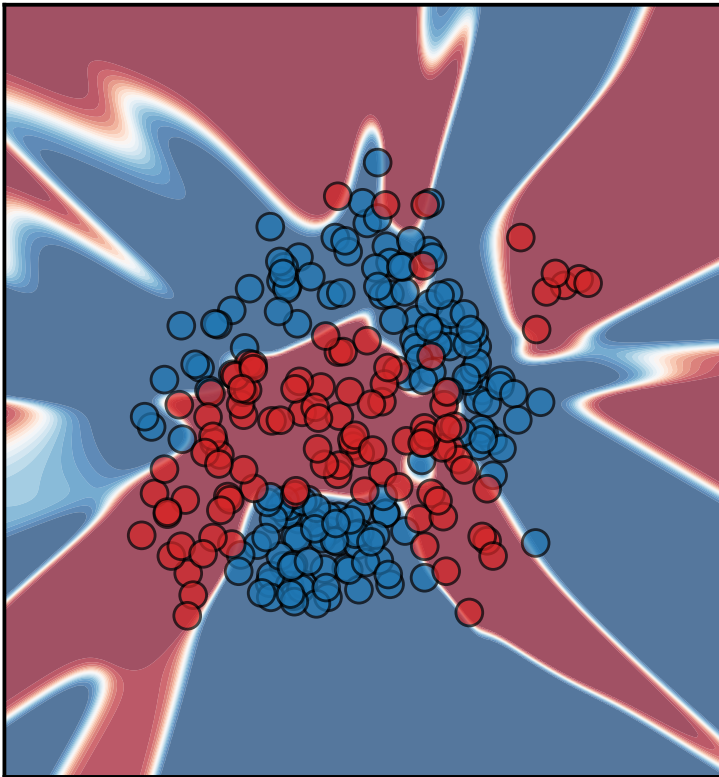
Uncertainty



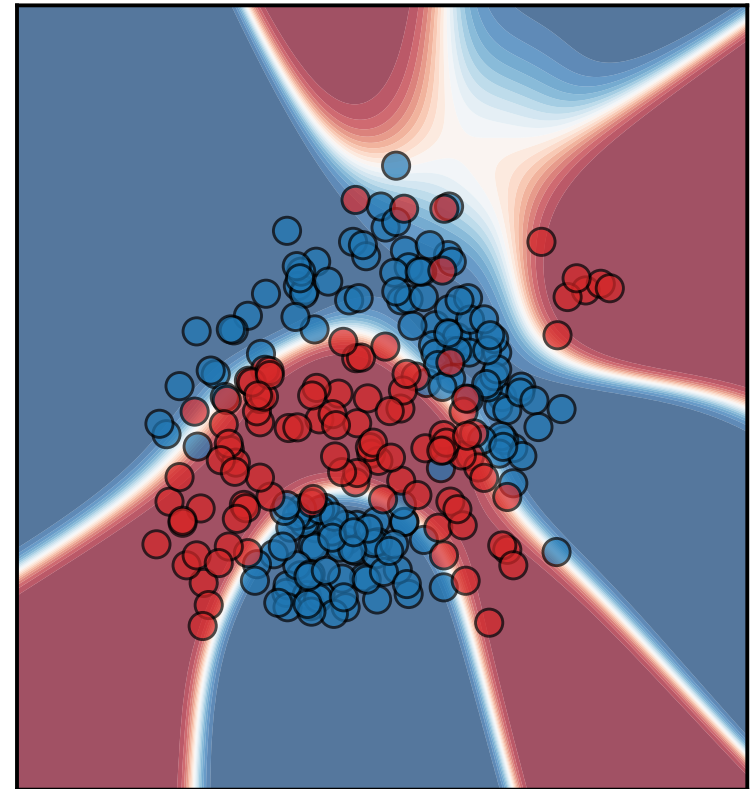
Kendall, Alex, Yarin Gal, and Roberto Cipolla. "Multi-task learning using uncertainty to weigh losses for scene geometry and semantics." *CVPR*. 2018.

# Reduce Overfitting

Standard DL



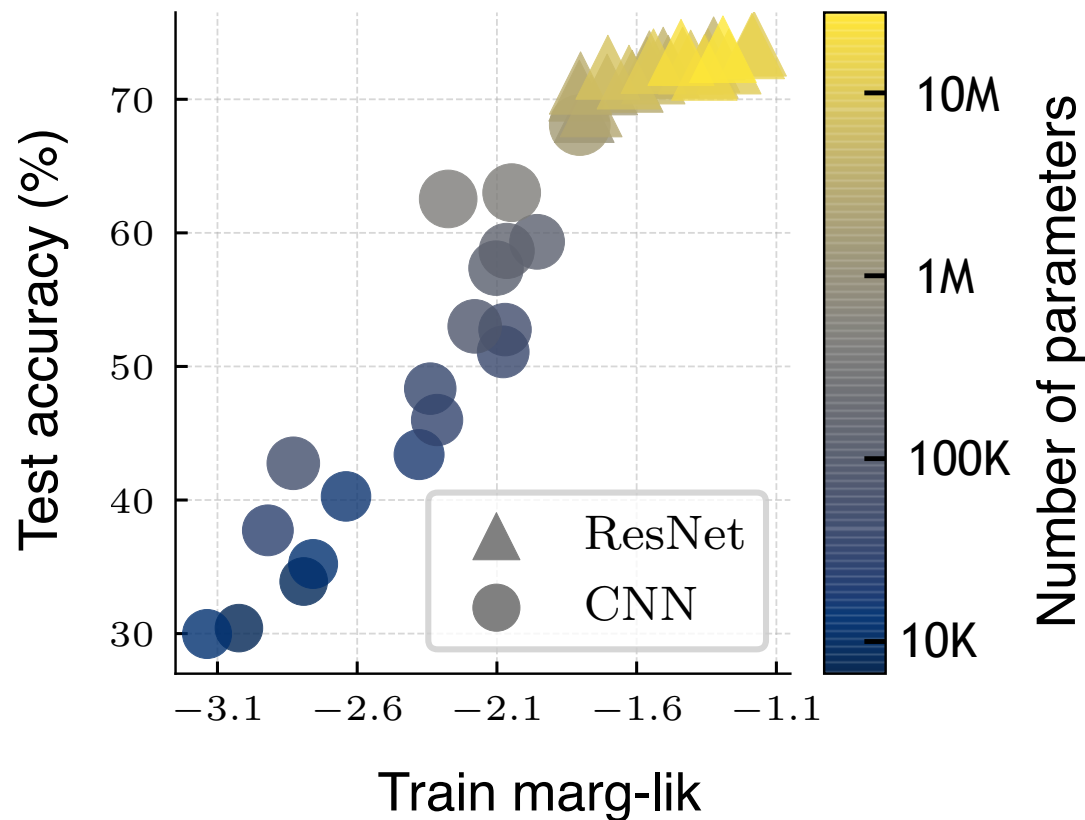
Bayesian DL



Left figure is cross-validation. Right figure is “Marginal Likelihood”.

# Model selection without test set

The “training marginal-likelihood” can be used to select deep-nets, *without* requiring the test set.



Test-accuracy correlates with train marg-lik.

Both increase as the model size is increased.

On CIFAR-100, around 50 models are shown.



# Bayesian learning

# Deep learning

Not scalable

Scalable

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{\int p(\mathcal{D}|\theta)p(\theta)d\theta}$$

$$\theta \leftarrow \theta - \rho H_{\theta}^{-1} \nabla_{\theta} \ell(\theta)$$

Bayesian Learning Rule:  $\lambda \leftarrow (1 - \rho)\lambda - \rho \nabla_{\mu} \mathbb{E}_q[\ell(\theta)]$

	Bayes	DL
Can handle large data and complex models?	✗	✓
Scalable training?	✗	✓
Can estimate uncertainty?	✓	✗
Can perform sequential / active /online / incremental learning?	✓	✗

# Bayesian Learning Rule

“Everything” from Bayes’ geometry  
New information as natural  
gradients



# The Bayesian Learning Rule

Mohammad Emtiyaz Khan  
RIKEN Center for AI Project  
Tokyo, Japan  
emtiyaz.khan@riken.jp

Håvard Rue  
CEMSE Division, KAUST  
Thuwal, Saudi Arabia  
haavard.rue@kaust.edu.sa

## Abstract

We show that many machine-learning algorithms are specific instances of a single algorithm called the *Bayesian learning rule*. The rule, derived from Bayesian principles, yields a wide-range of algorithms from fields such as optimization, deep learning, and graphical models. This includes classical algorithms such as ridge regression, Newton's method, and Kalman filter, as well as modern deep-learning algorithms such as stochastic-gradient descent, RMSprop, and Dropout. The key idea in deriving such algorithms is to approximate the posterior using candidate distributions estimated by using natural gradients. Different candidate distributions result in different algorithms and further approximations to natural gradients give rise to variants of those algorithms. Our work not only unifies, generalizes, and improves existing algorithms, but also helps us design new ones.

# Bayesian learning rule

See Table 1 in Khan and Rue, 2021

Learning Algorithm	Posterior Approx.	Natural-Gradient Approx.	Sec.
<b>Optimization Algorithms</b>			
Gradient Descent	Gaussian (fixed cov.)	Delta method	1.3
Newton's method	Gaussian	—“—	1.3
Multimodal optimization <sub>(New)</sub>	Mixture of Gaussians	—“—	3.2
<b>Deep-Learning Algorithms</b>			
Stochastic Gradient Descent	Gaussian (fixed cov.)	Delta method, stochastic approx.	4.1
RMSprop/Adam	Gaussian (diagonal cov.)	Delta method, stochastic approx., Hessian approx., square-root scaling, slow-moving scale vectors	4.2
Dropout	Mixture of Gaussians	Delta method, stochastic approx., responsibility approx.	4.3
STE	Bernoulli	Delta method, stochastic approx.	4.5
Online Gauss-Newton (OGN) <sub>(New)</sub>	Gaussian (diagonal cov.)	Gauss-Newton Hessian approx. in Adam & no square-root scaling	4.4
Variational OGN <sub>(New)</sub>	—“—	Remove delta method from OGN	4.4
BayesBiNN <sub>(New)</sub>	Bernoulli	Remove delta method from STE	4.5
<b>Approximate Bayesian Inference Algorithms</b>			
Conjugate Bayes	Exp-family	Set learning rate $\rho_t = 1$	5.1
Laplace's method	Gaussian	Delta method	4.4
Expectation-Maximization	Exp-Family + Gaussian	Delta method for the parameters	5.2
Stochastic VI (SVI)	Exp-family (mean-field)	Stochastic approx., local $\rho_t = 1$	5.3
VMP	—“—	$\rho_t = 1$ for all nodes	5.3
Non-Conjugate VMP	—“—	—“—	5.3
Non-Conjugate VI <sub>(New)</sub>	Mixture of Exp-family	None	5.4

# Principle of Trial-and-Error

Frequentist: Empirical Risk Minimization (ERM) or Maximum Likelihood Principle, etc.

$$\min_{\theta} \ell(\mathcal{D}, \theta) = \sum_{i=1}^N [y_i - f_{\theta}(x_i)]^2 + \gamma \theta^T \theta$$

Loss ↑  
Data ↑  
Model Params ↑  
Deep Network ↑

Deep Learning Algorithms:  $\theta \leftarrow \theta - \rho H_{\theta}^{-1} \nabla_{\theta} \ell(\theta)$

We will derive them as special instances of a rule exploiting information geometry of Bayes.

# Geometry of Exponential Family

We will exploit the geometry of “minimal” exp-family

Natural  
parameters

Sufficient  
Statistics

Expectation  
parameters

$$q(\theta) \propto \exp \left[ \lambda^\top T(\theta) \right]$$

$$\mu := \mathbb{E}_q[T(\theta)]$$

$$\begin{aligned} \mathcal{N}(\theta|m, S^{-1}) &\propto \exp \left[ -\frac{1}{2}(\theta - m)^\top S(\theta - m) \right] \\ &\propto \exp \left[ (Sm)^\top \theta + \text{Tr} \left( -\frac{S}{2} \theta \theta^\top \right) \right] \end{aligned}$$

Gaussian distribution

$$q(\theta) := \mathcal{N}(\theta|m, S^{-1})$$

Natural parameters

$$\lambda := \{Sm, -S/2\}$$

Expectation parameters

$$\mu := \{\mathbb{E}_q(\theta), \mathbb{E}_q(\theta \theta^\top)\}$$

1. Wainwright and Jordan, Graphical Models, Exp Fams, and Variational Inference Graphical models 2008

2. Malago et al., Towards the Geometry of Estimation of Distribution Algos based on Exp-Fam, FOGA, 2011 28



# The Bayesian Learning Rule

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

$\uparrow$   
 Posterior approximation (expo-family)

$\mathcal{H}(q)$   
 Entropy

**Bayesian Learning Rule** [1,2] (natural-gradient descent)

Natural and Expectation parameters of  $q$

$$\lambda \leftarrow \lambda - \rho \nabla_{\mu} \left\{ \mathbb{E}_q[\ell(\theta)] - \mathcal{H}(q) \right\}$$

$$\lambda \leftarrow (1 - \rho) \lambda - \rho \nabla_{\mu} \mathbb{E}_q[\ell(\theta)]$$

Old belief

New information = natural gradients

Exploiting posterior's information geometry to derive existing algorithms as special instances by approximating  $q$  and natural gradients.

1. Khan and Rue, The Bayesian Learning Rule, arXiv, <https://arxiv.org/abs/2107.04562>, 2021

2. Khan and Lin. "Conjugate-computation variational inference...." Alstats (2017).

# Warning!

- This natural gradient might be different from the one what we (often) encounter in machine learning for Maximum-Likelihood
  - In MLE, the loss is the negative log probability distribution

$$\min_{\theta} -\log q(\theta) \Rightarrow F(\theta)^{-1} \nabla \log q(\theta)$$

- Here,  $\theta$ , loss and distribution are two different entities, even possible unrelated

$$\min_q \mathbb{E}_q[\ell(\theta)] - \mathcal{H}(q) \Rightarrow F(\lambda)^{-1} \nabla_{\lambda} \mathbb{E}_q[\ell(\theta)]$$

# Gradient Descent from Bayesian Learning Rule

(Euclidean) gradients as natural  
gradients

# Bayesian learning rule:

Learning Algorithm	Posterior Approx.	Natural-Gradient Approx.	Sec.
<b>Optimization Algorithms</b>			
Gradient Descent	Gaussian (fixed cov.)	Delta method	1.3
Newton's method	Gaussian	—"—	1.3
Multimodal optimization <sub>(New)</sub>	Mixture of Gaussians	—"—	3.2
<b>Deep-Learning Algorithms</b>			
Stochastic Gradient Descent	Gaussian (fixed cov.)	Delta method, stochastic approx.	4.1
RMSprop/Adam	Gaussian (diagonal cov.)	Delta method, stochastic approx., Hessian approx., square-root scaling, slow-moving scale vectors	4.2
Dropout	Mixture of Gaussians	Delta method, stochastic approx., responsibility approx.	4.3
STE	Bernoulli	Delta method, stochastic approx.	4.5
Online Gauss-Newton (OGN) <sub>(New)</sub>	Gaussian (diagonal cov.)	Gauss-Newton Hessian approx. in Adam & no square-root scaling	4.4
Variational OGN <sub>(New)</sub>	—"—	Remove delta method from OGN	4.4
BayesBiNN <sub>(New)</sub>	Bernoulli	Remove delta method from STE	4.5
<b>Approximate Bayesian Inference Algorithms</b>			
Conjugate Bayes	Exp-family	Set learning rate $\rho_t = 1$	5.1
Laplace's method	Gaussian	Delta method	4.4
Expectation-Maximization	Exp-Family + Gaussian	Delta method for the parameters	5.2
Stochastic VI (SVI)	Exp-family (mean-field)	Stochastic approx., local $\rho_t = 1$	5.3
VMP	—"—	$\rho_t = 1$ for all nodes	5.3
Non-Conjugate VMP	—"—	—"—	5.3
Non-Conjugate VI <sub>(New)</sub>	Mixture of Exp-family	None	5.4

# Gradient Descent from BLR

$$\text{GD: } \theta \leftarrow \theta - \rho \nabla_{\theta} \ell(\theta)$$

$$\text{BLR: } m \leftarrow m - \rho \nabla_m \ell(m)$$

“Global” to “local”  
(the delta method)

$$\mathbb{E}_q[\ell(\theta)] \approx \ell(m)$$

$$m \leftarrow m - \rho \nabla_m \mathbb{E}_q[\ell(\theta)]$$

$$\lambda \leftarrow \lambda - \rho \nabla_{\mu} (\mathbb{E}_q[\ell(\theta)] - \mathcal{H}(q))$$

Derived by choosing **Gaussian with fixed covariance**

Gaussian distribution  $q(\theta) := \mathcal{N}(m, 1)$

Natural parameters  $\lambda := m$

Expectation parameters  $\mu := \mathbb{E}_q[\theta] = m$

Entropy  $\mathcal{H}(q) := \log(2\pi)/2$

# Bayesian learning rule:

Put the expectation (Bayes) back in and use the Bayesian averaging.

Learning Algorithm	Posterior Approx.	Natural-Gradient Approx.	Sec.
<b>Optimization Algorithms</b>			
Gradient Descent	Gaussian (fixed cov.)	Delta method	1.3
Newton's method	Gaussian	—"—"	1.3
Multimodal optimization <small>(New)</small>	Mixture of Gaussians	—"—"	3.2
<b>Deep-Learning Algorithms</b>			
Stochastic Gradient Descent	Gaussian (fixed cov.)	Delta method, stochastic approx.	4.1
RMSprop/Adam	Gaussian (diagonal cov.)	Delta method, stochastic approx., Hessian approx., square-root scaling, slow-moving scale vectors	4.2
Dropout	Mixture of Gaussians	Delta method, stochastic approx., responsibility approx.	4.3
STE	Bernoulli	Delta method, stochastic approx.	4.5
Online Gauss-Newton (OGN) <small>(New)</small>	Gaussian (diagonal cov.)	Gauss-Newton Hessian approx. in Adam & no square-root scaling	4.4
Variational OGN <small>(New)</small>	—"—"	Remove delta method from OGN	4.4
BayesBiNN <small>(New)</small>	Bernoulli	Remove delta method from STE	4.5
<b>Approximate Bayesian Inference Algorithms</b>			
Conjugate Bayes	Exp-family	Set learning rate $\rho_t = 1$	5.1
Laplace's method	Gaussian	Delta method	4.4
Expectation-Maximization	Exp-Family + Gaussian	Delta method for the parameters	5.2
Stochastic VI (SVI)	Exp-family (mean-field)	Stochastic approx., local $\rho_t = 1$	5.3
VMP	—"—"	$\rho_t = 1$ for all nodes	5.3
Non-Conjugate VMP	—"—"	—"—"	5.3
Non-Conjugate VI <small>(New)</small>	Mixture of Exp-family	None	5.4

1. Khan, et al. "Fast and scalable Bayesian deep learning by weight-perturbation in Adam." *ICML* (2018).
2. Osawa et al. "Practical Deep Learning with Bayesian Principles." *NeurIPS* (2019).
3. Lin et al. "Handling the positive-definite constraints in the BLR." *ICML* (2020).

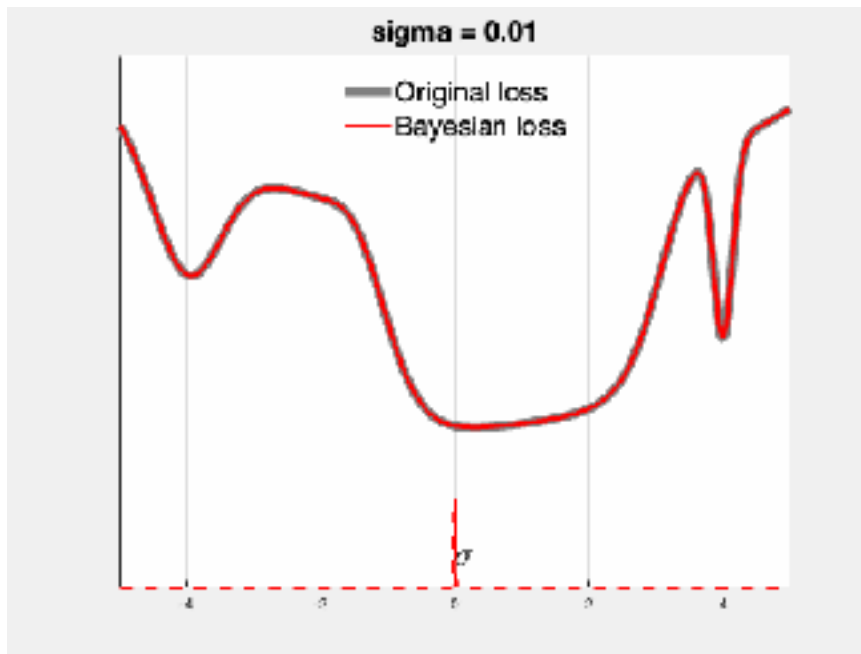
# Why use Bayesian averaging?

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

$\uparrow$   
 Gaussian approximation

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

Entropy



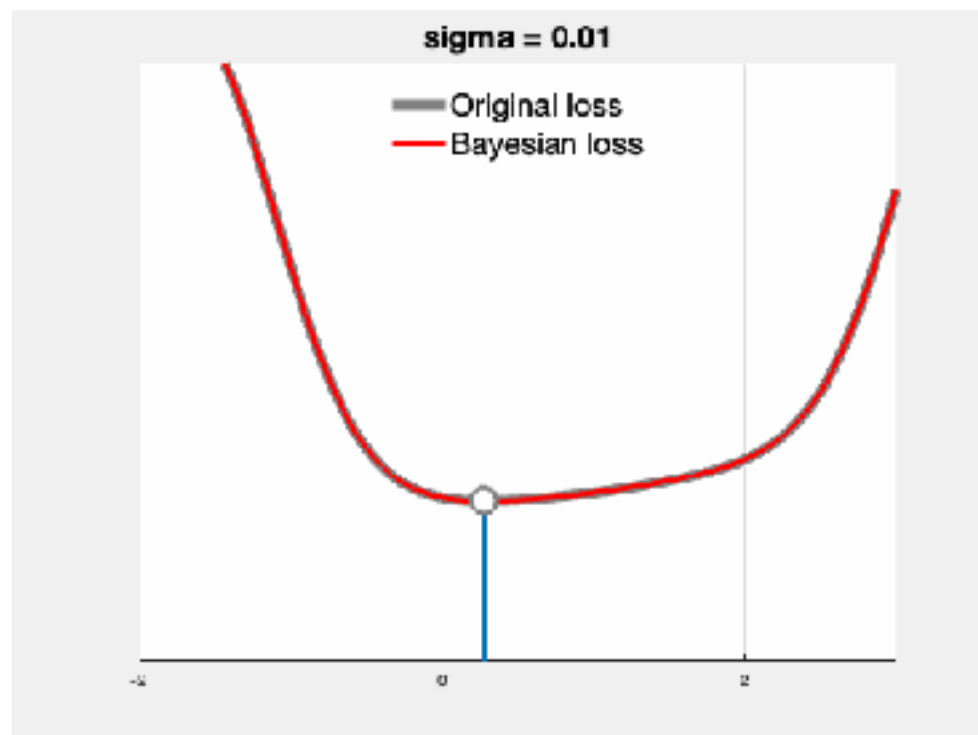
- First term “smooths” the loss to favor “flatter” regions [1]
- Second term figures out how much to smooth (find  $\sigma$ )
- Similar mechanisms are used in DL [2], RL, search, robust optimization.

# Bayes Prefers Flatter directions

$$\text{GD: } \theta \leftarrow \theta - \rho \nabla_{\theta} \ell(\theta) \quad \implies \nabla_{\theta} \ell(\theta_*) = 0$$

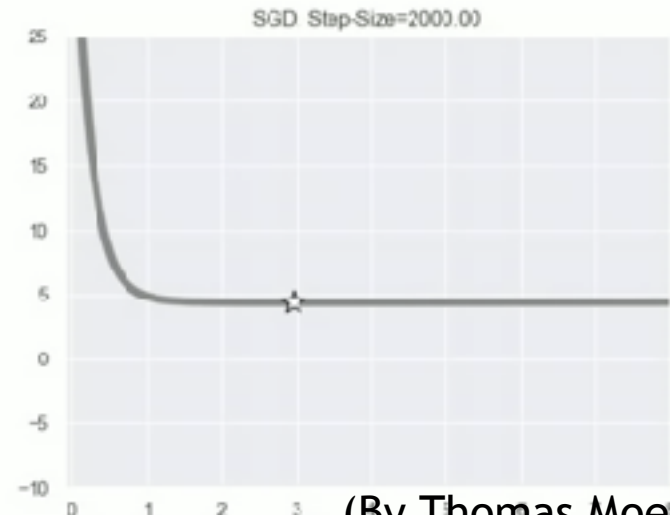
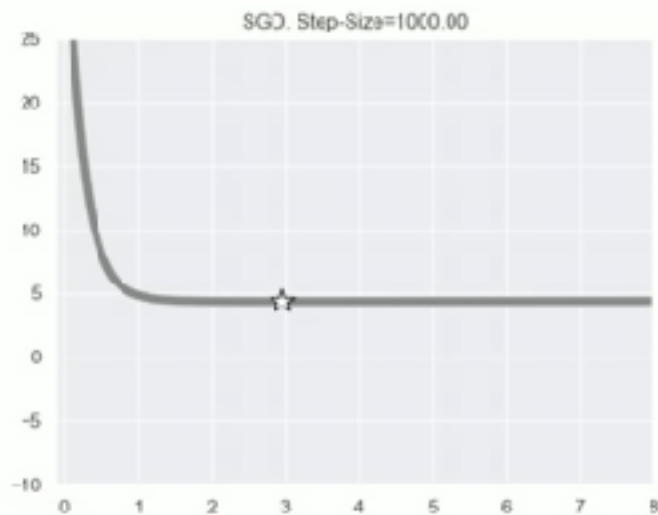
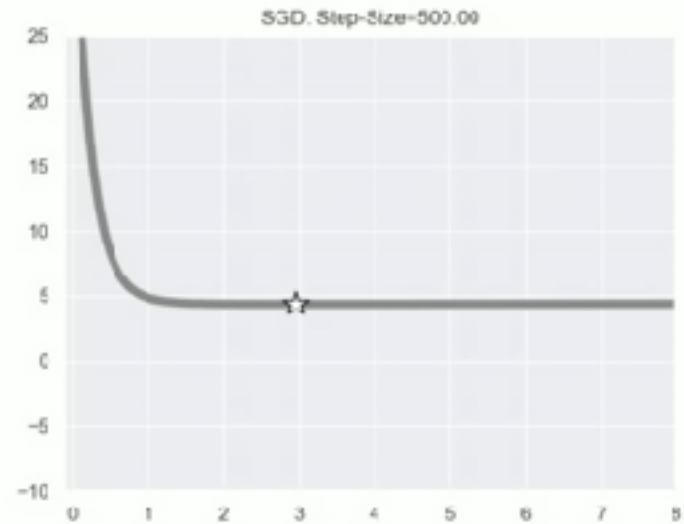
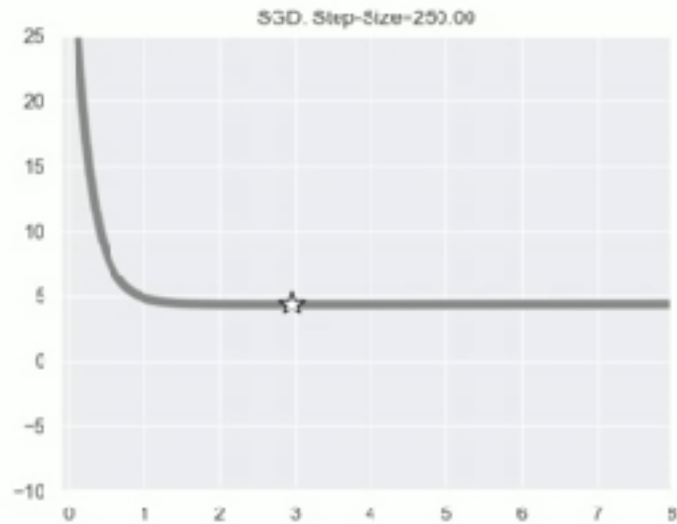
$$\text{BLR: } m \leftarrow m - \rho \nabla_m \mathbb{E}_q[\ell(\theta)] \quad \implies \nabla_m \mathbb{E}_{q^*}[\ell(\theta)] = 0$$

Bayesian solution injects “noise” which has a similar regularization effect to noise in Stochastic GD. It prefers “flatter” directions.

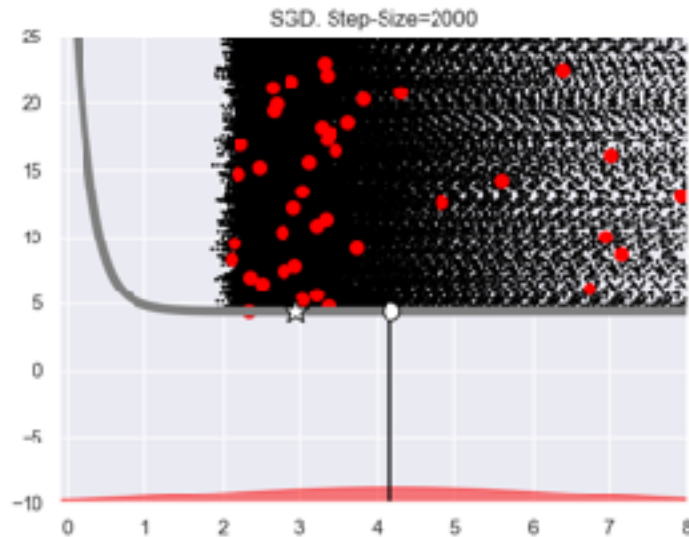
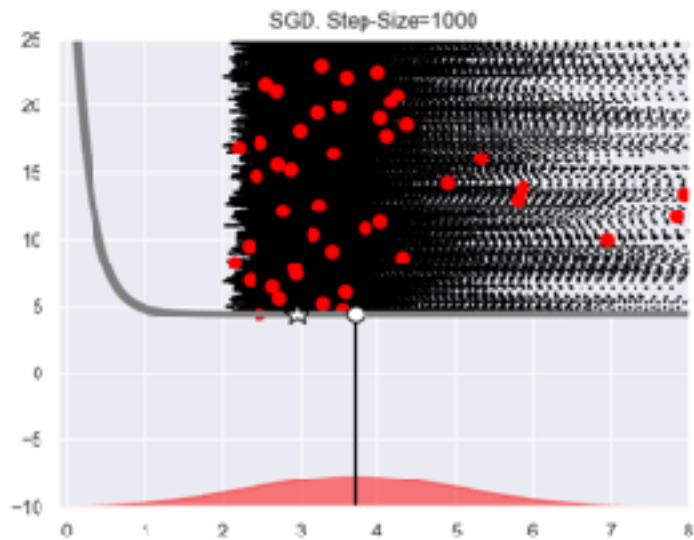
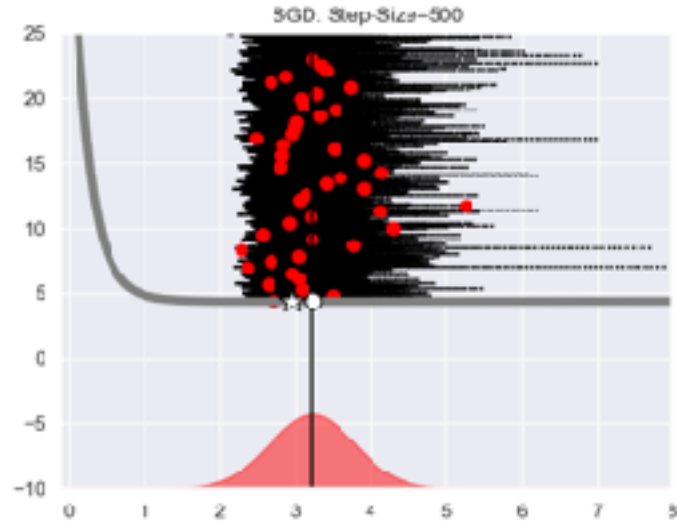
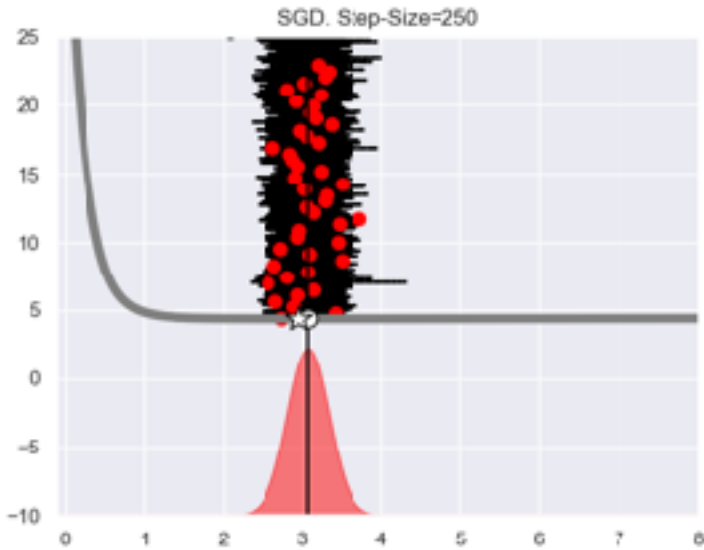




# SGD: Implicit Regularization



# SGD: Implicit Regularization



# Bayes: Explicit Regularization

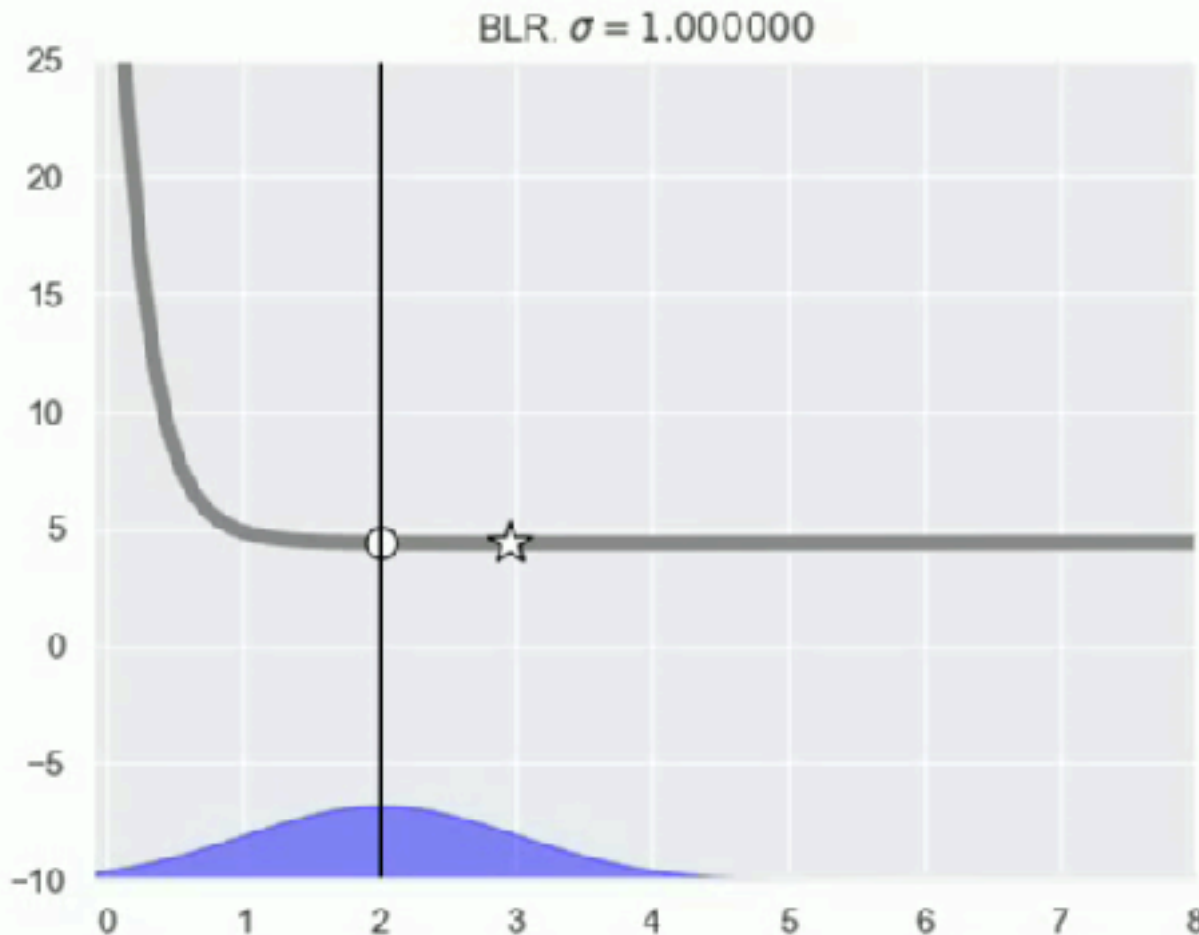
Estimating Gaussian posteriors where the variance is fixed, and only the mean is estimated

$$\mathbb{E}_{q_*}[\nabla_{\theta} \ell(\theta)] = 0$$

By increasing the variance, we can move the mode arbitrarily far.

Bayesian “noise” has a similar regularization to the SGD noise.

It prefers “flatter” directions.



# **Newton's method from Bayesian Learning Rule**

(Gradient, Hessian) as natural  
gradients

# Newton's Method from BLR

Newton's method:  $\theta \leftarrow \theta - H_\theta^{-1} [\nabla_\theta \ell(\theta)]$

$$Sm \leftarrow (1 - \rho)Sm - \rho \nabla_{\mathbb{E}_q(\theta)} \mathbb{E}_q[\ell(\theta)]$$

$$-\frac{1}{2}S \leftarrow (1 - \rho)S - \rho \frac{1}{2} S \nabla_{\mathbb{E}_q(\theta)} \nabla_{\mathbb{E}_q(\theta)} \mathbb{E}_q[\ell(\theta)]$$

$$\lambda \leftarrow (1 - \rho) \nabla_{\mu} \mathbb{E}_q[\nabla_{\mu} \ell(\theta)] - \rho \nabla_{\mu} \mathbb{E}_q[\nabla_{\mu} \ell(\theta)] \quad -\nabla_{\mu} \mathcal{H}(q) = \lambda$$

Derived by choosing a **multivariate Gaussian**

Gaussian distribution  $q(\theta) := \mathcal{N}(\theta|m, S^{-1})$

Natural parameters  $\lambda := \{Sm, -S/2\}$

Expectation parameters  $\mu := \{\mathbb{E}_q(\theta), \mathbb{E}_q(\theta\theta^\top)\}$

# Newton's Method from BLR

Newton's method:  $\theta \leftarrow \theta - H_{\theta}^{-1} [\nabla_{\theta} \ell(\theta)]$

Set  $\rho = 1$  to get  $m \leftarrow m - H_m^{-1} [\nabla_m \ell(m)]$

$$m \leftarrow m - \rho S^{-1} \nabla_m \ell(m)$$

$$S \leftarrow (1 - \rho)S + \rho H_m$$

Delta Method

$$\mathbb{E}_q[\ell(\theta)] \approx \ell(m)$$

Express in terms of gradient and Hessian of loss:

$$\nabla_{\mathbb{E}_q(\theta)} \mathbb{E}_q[\ell(\theta)] = \mathbb{E}_q[\nabla_{\theta} \ell(\theta)] - 2\mathbb{E}_q[H_{\theta}]m$$

$$\nabla_{\mathbb{E}_q(\theta\theta^{\top})} \mathbb{E}_q[\ell(\theta)] = \mathbb{E}_q[H_{\theta}]$$

$$Sm \leftarrow (1 - \rho)Sm - \rho \nabla_{\mathbb{E}_q(\theta)} \mathbb{E}_q[\ell(\theta)]$$

$$S \leftarrow (1 - \rho)S - \rho 2 \nabla_{\mathbb{E}_q(\theta\theta^{\top})} \mathbb{E}_q[\ell(\theta)]$$

# RMSprop/Adam from BLR

RMSprop

$$s \leftarrow (1 - \rho)s + \rho[\hat{\nabla} \ell(\theta)]^2$$

$$\theta \leftarrow \theta - \alpha(\sqrt{s} + \delta)^{-1} \hat{\nabla} \ell(\theta)$$

BLR for Gaussian approx

$$S \leftarrow (1 - \rho)S + \rho(H_\theta)$$

$$m \leftarrow m - \alpha S^{-1} \nabla_\theta \ell(\theta)$$

To get RMSprop, make the following choices

- Restrict covariance to be diagonal
- Replace Hessian by square of gradients
- Add square root for scaling vector

For Adam, use a Heavy-ball term with KL divergence as momentum (Appendix E in [1])

# Practical DL with Bayes

## RMSprop

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

$$s \leftarrow (1 - \rho)s + \rho g^2$$

$$\theta \leftarrow \theta - \alpha(\sqrt{s} + \delta)^{-1} g$$

## BLR variant called VOGN

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

$$s \leftarrow (1 - \rho)s + \rho(\sum_i g_i^2)$$

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

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

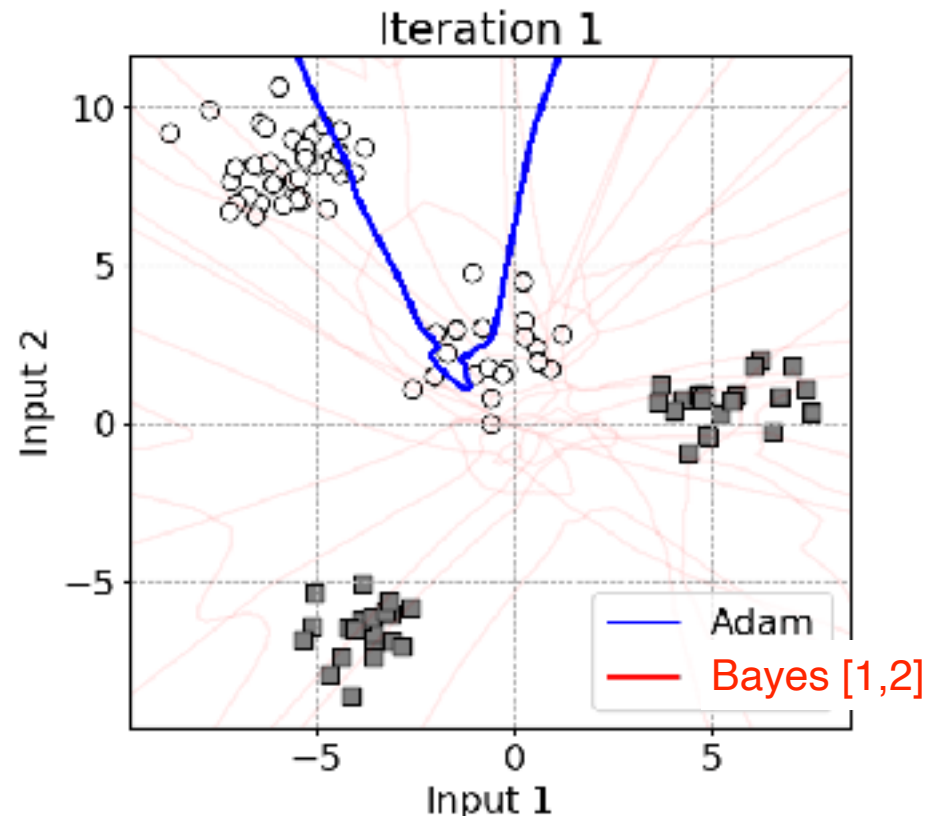
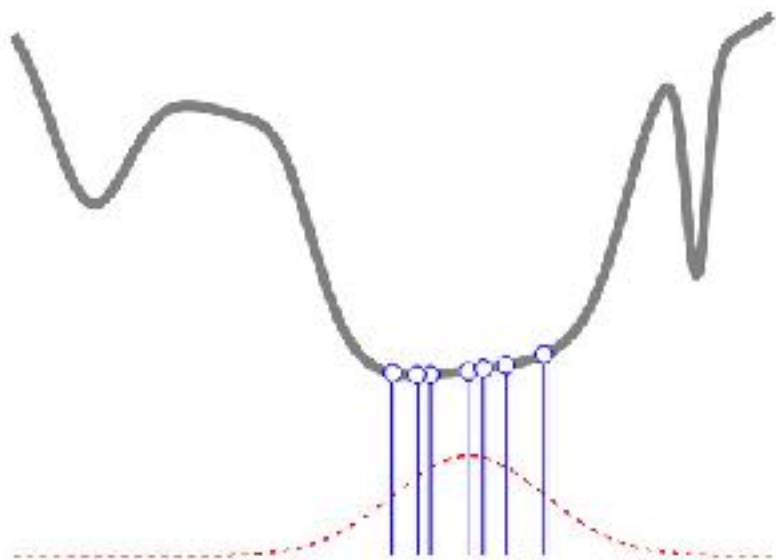
Available at <https://github.com/team-approx-bayes/dl-with-bayes>

1. Khan, et al. "Fast and scalable Bayesian deep learning by weight-perturbation in Adam." *ICML* (2018).
2. Osawa et al. "Practical Deep Learning with Bayesian Principles." *NeurIPS* (2019).
3. Lin et al. "Handling the positive-definite constraints in the BLR." *ICML* (2020).



# Why use Bayesian averaging?

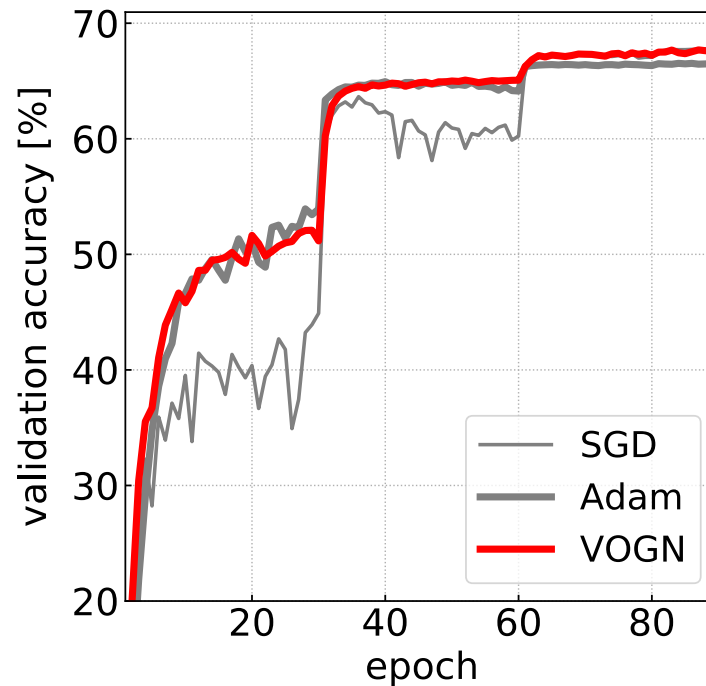
Choose an “ensemble” of almost equally good models (similar to sampling in SGD trajectories)



1. Khan, et al. "Fast and scalable Bayesian deep learning by weight-perturbation in Adam." *ICML* (2018).
2. Osawa et al. "Practical Deep Learning with Bayesian Principles." *NeurIPS* (2019).

# Uncertainty of Deep Nets

VOGN: A modification of Adam with similar performance on ImageNet, but better uncertainty



Code available at <https://github.com/team-approx-bayes/dl-with-bayes>

1. Khan, et al. "Fast and scalable Bayesian deep learning by weight-perturbation in Adam." *ICML* (2018).
2. Osawa et al. "Practical Deep Learning with Bayesian Principles." *NeurIPS* (2019).

# BLR variant [3] got 1st prize in NeurIPS 2021 Approximate Inference Challenge

Watch **Thomas Moellenhoff's** talk at <https://www.youtube.com/watch?v=LQInIN5EU7E>.

## Mixture-of-Gaussian Posteriors with an Improved Bayesian Learning Rule

Thomas Möllenhoff<sup>1</sup>, Yuesong Shen<sup>2</sup>, Gian Maria Marconi<sup>1</sup>  
Peter Nickl<sup>1</sup>, Mohammad Emtiyaz Khan<sup>1</sup>

  
**1** Approximate Bayesian Inference Team  
RIKEN Center for AI Project, Tokyo, Japan

  
**2** Computer Vision Group  
Technical University of Munich, Germany





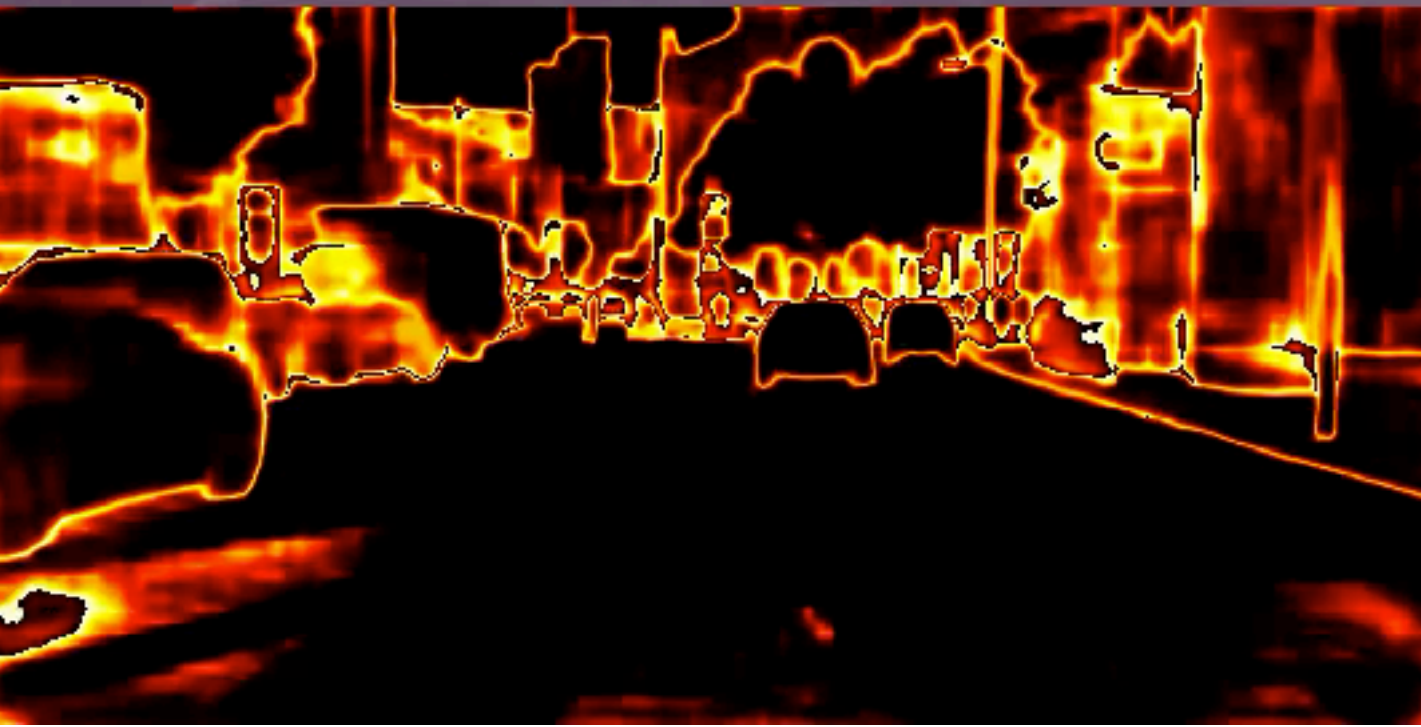


Dec 14th, 2021 — NeurIPS Workshop on Bayesian Deep Learning

1. Khan, et al. "Fast and scalable Bayesian deep learning by weight-perturbation in Adam." *ICML* (2018).
2. Osawa et al. "Practical Deep Learning with Bayesian Principles." NeurIPS (2019).
3. Lin et al. "Handling the positive-definite constraints in the BLR." *ICML* (2020).



Image  
Segmentation



Uncertainty  
(entropy of  
class probs)

# Summary

- Gradient descent is derived using a Gaussian with fixed covariance, and estimating the mean
- Newton's method is derived using multivariate Gaussian
- RMSprop is derived using diagonal covariance
- Adam is derived by adding heavy-ball momentum term
- For “ensemble of Newton”, use Mixture of Gaussians [1]
- To derive DL algorithms, we need to use the Delta method (a local approximation)  $\mathbb{E}_q[\ell(\theta)] \approx \ell(m)$
- Then, **to improve DL algorithms, we just need to add some “global” touch by relaxing the local approximation**

1. Lin, Wu, Mohammad Emtiyaz Khan, and Mark Schmidt. "Fast and Simple Natural-Gradient Variational Inference with Mixture of Exponential-family Approximations." *ICML* (2019).

# Bayes' Rule from Bayesian Learning Rule

“Messages” as natural gradients

# Bayesian Inference as Optimization

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{\int p(\mathcal{D}|\theta)p(\theta)d\theta}$$

$$\ell(\theta) := -\log p(\mathcal{D}|\theta)p(\theta)$$

$$= \arg \min_{q \in \mathcal{P}} \mathbb{E}_{q(\theta)}[\ell(\theta)] - \mathcal{H}(q)$$

All distribution      Distribution      Entropy

$$= \mathbb{E}_q[\ell(\theta)] + \mathbb{E}_q[\log q(\theta)] = \mathbb{E}_q \left[ \log \frac{q(\theta)}{e^{-\ell(\theta)}} \right]$$
$$\implies q_*(\theta) \propto e^{-\ell(\theta)} \propto p(\mathcal{D}|\theta)p(\theta) \propto p(\theta|\mathcal{D})$$

Good news: This holds for a generic loss function!

# Bayesian Inference from BLR

Ex: Linear model, Kalman filters, HMM, etc.

$$\ell(\theta) := -\log p(\mathcal{D}|\theta)p(\theta) = -\lambda_{\mathcal{D}}^{\top} T(\theta) \leftarrow \begin{array}{l} \text{Sufficient} \\ \text{statistics of } q \end{array}$$

$$\begin{aligned} \ell(\theta) &:= (y - X\theta)^{\top} (y - X\theta) + \gamma\theta^{\top} \theta \\ &= -2\theta^{\top} (X^{\top} y) + \text{Tr} [\theta\theta^{\top} (X^{\top} X + \gamma I)] + \text{cnst} \end{aligned}$$

$$\implies \mathbb{E}_q[\ell(\theta)] = -\lambda_{\mathcal{D}} \mu \implies \nabla_{\mu} \mathbb{E}_q[\ell(\theta)] = -\lambda_{\mathcal{D}}$$

$$\lambda \leftarrow \lambda - \rho (\nabla_{\mu} \lambda \mathbb{E}_q[\ell(\theta)] + \lambda \theta(q)) \implies \lambda_* = \lambda_{\mathcal{D}}$$

Forward-backward, SVI, Variational message passing etc. are special cases of the BLR Messages



# The Bayesian “Principle”

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

All distribution  $\uparrow$  Distribution  $\uparrow$  Entropy

$$\ell(\theta) := -\log p(\mathcal{D}|\theta)p(\theta)$$

Restrict the set of distributions (change P to Q)

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

Often we call this variational Inference, but it's not just a method, rather a general Bayesian principle.

Similar to IGO [1] but the entropy is essential!

Our use of natural-gradients here is not a matter of choice. In fact, natural-gradients are inherently present in *all solutions of the Bayesian objective* in Eq. 2. For example, a solution of Eq. 2 or equivalently a fixed point of Eq. 3, satisfies the following,

$$\nabla_{\mu} \mathbb{E}_{q_*}[\tilde{\ell}(\theta)] = \nabla_{\mu} \mathcal{H}(q_*), \text{ which implies } \tilde{\nabla}_{\lambda} \mathbb{E}_{q_*}[-\tilde{\ell}(\theta)] = \lambda_*, \quad (5)$$

for candidates with constant base-measure. This is obtained by setting the gradient of Eq. 2 to 0, then noting that  $\nabla_{\mu} \mathcal{H}(q) = -\lambda$  (App. B), and then interchanging  $\nabla_{\mu}$  by  $\tilde{\nabla}_{\lambda}$  (because of Eq. 4). In other words, natural parameter of the best  $q_*(\theta)$  is equal to the natural gradient of the expected negative-loss. The importance of natural-gradients is entirely missed in the Bayesian/variational inference literature, including textbooks, reviews, tutorials on this topic [Bishop, 2006, Murphy, 2012, Blei et al., 2017, Zhang et al., 2018a] where natural-gradients are often put in a special category.

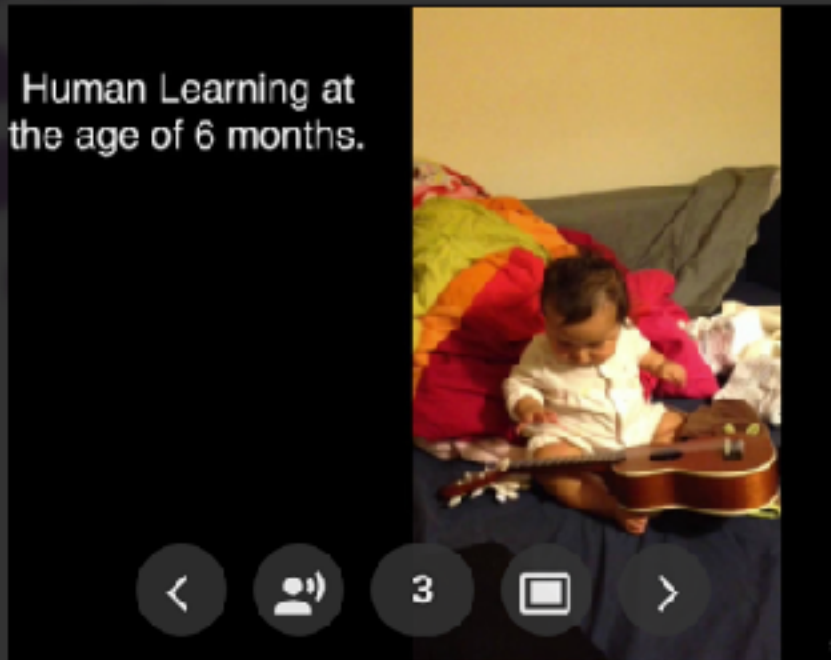
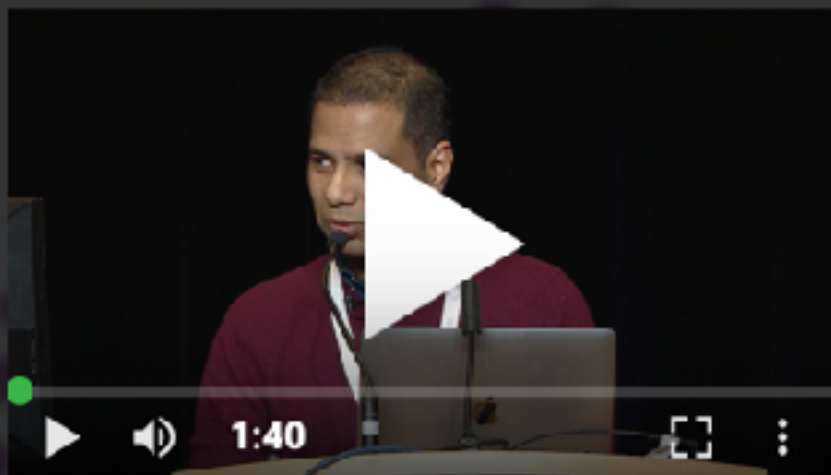
We will show that natural gradients retrieve essential higher-order information about the loss landscape which are then assigned to appropriate natural parameters using Eq. 5. The information-matching is due to the presence of the entropy term there, which is an important quantity for the optimality of Bayes in general [Jaynes, 1982, Zellner, 1988, Littlestone and Warmuth, 1994, Vovk, 1990], and which is generally absent in non-Bayesian formulations (Eq. 1). The entropy term in general leads to exponential-weighting in Bayes' rule. In our context, it gives rise to natural-gradients and, as we will soon see, automatically determines the complexity of the derived algorithm through the complexity of the class of distributions  $\mathcal{Q}$ , yielding a principled way to develop new algorithms.

Overall, our work demonstrates the importance of natural-gradients and information geometry for algorithm design in ML. This is similar in spirit to Information Geometric Optimization [Ollivier et al., 2017], which focuses on the optimization of black-box, deterministic functions. In contrast, we derive generic learning algorithms by using the same Bayesian principles. The BLR we use is a generalization of the method proposed in Khan and Lin [2017], Khan and Nielsen [2018] specifically for approximate Bayesian inference. Here, we establish it as a general learning rule to derive many old and new learning algorithms, which include both Bayesian and non-Bayesian ones, way beyond its original proposal. We do not claim that these successful algorithms work well because they are derived from the BLR. Rather, we use the BLR to simply unravels the inherent Bayesian nature of these "good" algorithms. In this sense, the BLR can be seen as a variant of Bayes' rule, useful for generic algorithm design.

# Principles of “good” algorithms?

- Information Geometry of Bayes
  - To unify/generalize/improve learning-algorithms
  - Optimize for “posterior approximations”
- Bayesian Learning rule (BLR)
  - Derive many algorithms from optimization, deep learning, and Bayesian inference
- Natural Gradients are Everywhere!

# NeurIPS 2019 Tutorial



## Deep Learning with Bayesian Principles

by **Mohammad Emtiyaz Khan** · Dec 9, 2019

#NeurIPS 2019

Follow

Views 151 807

Presentations 263

Followers 200

Latest

Popular

...



From System 1 Deep Learning to System 2 Deep Learning

by [Yoshua Bengio](#)

17,953 views · Dec 11, 2019



NeurIPS Workshop on Machine Learning for Creativity and Design...

by [Aaron Hertzmann](#) [Adam Roberts](#) ...

9,654 views · Dec 14, 2019



Deep Learning with Bayesian Principles

by [Mohammad Emtiyaz Khan](#)

4,084 views · Dec 5, 2019



Efficient Processing of Deep Neural Network: from Algorithms to...

by [Wenyuan Szeg](#)

7,162 views · Dec 9, 2019

# What's Next

- Bayesian “Duality” Principle
  - The BLR unravels a duality perspective of good algorithms
  - Unifies many results from many fields
    - convex duality, Kernel methods, Bayesian nonparametric methods, Deep Learning, Robust statistics, and Information Geometry
  - Helps to “solve” the Adaptation problem
- My talk on Monday will show two examples of this principle
  - Robust deep learning as “convex relaxation” of Bayes
  - Principle of Adaptive learning through K-priors

# 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**

Co-PI  
(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 Augier**  
Research Scientist



**Hugo Monzón  
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



**Jaeyoon Kim**  
Postdoc



**Alexandre Piché**  
Remote Collaborator  
MILA



**All Unlu**  
Intern, Okinawa  
Institute of Science



**Geoffrey Wolfer**  
Postdoc



**Wu Lin**  
PhD Student  
University of British  
Columbia



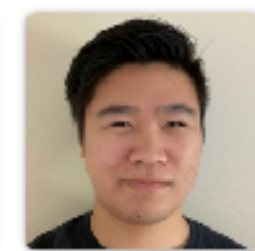
**Peter Nickl**  
Research Assistant



**Dharmesh Tailor**  
Remote Collaborator  
University of  
Amsterdam



**Ang Mingliang**  
Remote Collaborator  
National University of  
Singapore



**Kenneth Chen**  
Intern, Okinawa  
Institute of Science  
and Technology