



Bayesian Principles for Learning-Machines

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



Al that learn like humans

Quickly adapt to learn new skills, throughout their lives

Human Learning at the age of 6 months.



Human Learning at the age of 6 months.



Human Learning at the age of 6 months.



Converged at the age of 12 months



Converged at the age of 12 months



Converged at the age of 12 months



Transfer skills at the age of 14 months



Transfer skills at the age of 14 months



Transfer skills at the age of 14 months



Failure of AI in "dynamic" setting

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



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

Failure of AI in "dynamic" setting

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



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

July 14, 2021



Yann LeCun @ylecun · 7h

So many exciting new frontiers in ML, it's hard to give a short list, particularly in new application areas (e.g. in the physical and biological sciences).

But the Big Question is:

"How could machines learn as efficiently as humans and animals?" This requires new paradigms.

Towards a new learning paradigm, based on Bayesian principles

Human learning

Life-long learning from small chunks of data in a non-stationary world

Deep learning

Bulk learning from a large amount of data in a stationary world

1. Parisi, German I., et al. "Continual lifelong learning with neural networks: A review." Neural Networks (2019)

Human learning

Life-long learning from small chunks of data in a non-stationary world

Deep learning

Bulk learning from a large amount of data in a stationary world

Our current research focuses on reducing this gap!

1. Parisi, German I., et al. "Continual lifelong learning with neural networks: A review." Neural Networks (2019)

Bayesian Principles

Human learning

Life-long learning from small chunks of data in a non-stationary world Deep learning

Bulk learning from a large amount of data in a stationary world

Our current research focuses on reducing this gap!

Parisi, German I., et al. "Continual lifelong learning with neural networks: A review." *Neural Networks* (2019)
 Geisler, W. S., and Randy L. D. "Bayesian natural selection and the evolution of perceptual systems." *Philosophical Transactions of the Royal Society of London. Biological Sciences* (2002)

Bayesian Principles

Human learning

Life-long learning from small chunks of data in a non-stationary world Deep learning

ur research

Bulk learning from a large amount of data in a stationary world

Our current research focuses on reducing this gap!

Parisi, German I., et al. "Continual lifelong learning with neural networks: A review." *Neural Networks* (2019)
 Geisler, W. S., and Randy L. D. "Bayesian natural selection and the evolution of perceptual systems." *Philosophical Transactions of the Royal Society of London. Biological Sciences* (2002)

- Bayesian principles as a general principle
 - To unify/generalize/improve learning-algorithms
 - By computing "posterior approximations"

- Bayesian principles as a general principle
 - To unify/generalize/improve learning-algorithms
 - By computing "posterior approximations"
- Bayesian Learning rule (BLR)
 - Derive many existing algorithms
 - Deep Learning (SGD, RMSprop, Adam)
 - Design new algorithms for uncertainty in DL

- Bayesian principles as a general principle
 - To unify/generalize/improve learning-algorithms
 - By computing "posterior approximations"
- Bayesian Learning rule (BLR)
 - Derive many existing algorithms
 - Deep Learning (SGD, RMSprop, Adam)
 - Design new algorithms for uncertainty in DL
- Dual perspective of BLR for life-long learning

- Bayesian principles as a general principle
 - To unify/generalize/improve learning-algorithms
 - By computing "posterior approximations"
- Bayesian Learning rule (BLR)
 - Derive many existing algorithms
 - Deep Learning (SGD, RMSprop, Adam)
 - Design new algorithms for uncertainty in DL
- Dual perspective of BLR for life-long learning
- Impact: Everything with the same principle



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.

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

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

$$\min_{\substack{\theta \text{ Loss } \\ \text{ Data }}} \ell(\mathcal{D}, \theta)$$

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

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

$$\sup_{\substack{i=1 \\ \mathsf{Deep} \\ \mathsf{Network}}} \int_{\mathsf{Deep}} \frac{1}{\mathsf{Network}}$$

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

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

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

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



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.

Bayes Objective

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

Sandard Deviation: 0.00 Particular deviation: 0.00 Parti

Instead of the original loss, optimize a different (smoothed) one (a very popular idea now for DL theory [4]).

A common idea in Inference, optimization, online learning, Reinforcement learning

Zellner, A. "Optimal information processing and Bayes's theorem." *The American Statistician* (1988)
 Many other: Bissiri, et al. (2016), Shawe-Taylor and Williamson (1997), Cesa-Bianchi and Lugosi (2006)
 Huszar's blog, Evolution Strategies, Variational Optimisation and Natural ES (2017)
 Smith et al., On the Origin of Implicit Regularization in Stochastic Gradient Descent, ICLR, 2021

Bayes Objective

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

Sandard Deviation: 0.00 Particular deviation: 0.00 Parti

Instead of the original loss, optimize a different (smoothed) one (a very popular idea now for DL theory [4]).

A common idea in Inference, optimization, online learning, Reinforcement learning

Zellner, A. "Optimal information processing and Bayes's theorem." *The American Statistician* (1988)
 Many other: Bissiri, et al. (2016), Shawe-Taylor and Williamson (1997), Cesa-Bianchi and Lugosi (2006)
 Huszar's blog, Evolution Strategies, Variational Optimisation and Natural ES (2017)
 Smith et al., On the Origin of Implicit Regularization in Stochastic Gradient Descent, ICLR, 2021

Step 1: Choose an approximation (mix-exp-family)

Natural parameters Sufficient statistics $q(\theta) \propto \exp\left[\lambda^{\top} T(\theta)\right]$

Step 1: Choose an approximation (mix-exp-family)

Natural parameters Sufficient statistics $q(\theta) \propto \exp\left[\lambda^{\top} T(\theta)\right]$ $\mathcal{N}(\theta|m, S^{-1}) \propto \exp\left[-\frac{1}{2}(\theta - m)^{\top} S(\theta - m)\right]$

Step 1: Choose an approximation (mix-exp-family)

Natural parameters Sufficient statistics $q(\theta) \propto \exp\left[\lambda^{\top} T(\theta)\right]$ $\mathcal{N}(\theta|m, S^{-1}) \propto \exp\left[-\frac{1}{2}(\theta - m)^{\top} S(\theta - m)\right]$ $\propto \exp\left[(Sm)^{\top} \theta + \operatorname{Tr}\left(-\frac{S}{2}\theta\theta^{\top}\right)\right]$

Step 1: Choose an approximation (mix-exp-family)

Natural parameters Sufficient statistics Expectation parameters

$$q(\theta) \propto \exp \left[\lambda^{\top} T(\theta)\right]$$
 \downarrow
 $\mu := \mathbb{E}_q[T(\theta)]$
 $\mathcal{N}(\theta|m, S^{-1}) \propto \exp \left[-\frac{1}{2}(\theta - m)^{\top} S(\theta - m)\right]$
 $\propto \exp \left[(Sm)^{\top} \theta + \operatorname{Tr} \left(-\frac{S}{2}\theta\theta^{\top}\right)\right]$

Step 1: Choose an approximation (mix-exp-family)

Natural parameters Sufficient statistics Expectation parameters

$$q(\theta) \propto \exp \left[\lambda^{\top} T(\theta)\right]$$
 \downarrow
 $\mu := \mathbb{E}_q[T(\theta)]$
 $\mathcal{N}(\theta|m, S^{-1}) \propto \exp \left[-\frac{1}{2}(\theta - m)^{\top} S(\theta - m)\right]$
 $\propto \exp \left[(Sm)^{\top} \theta + \operatorname{Tr}\left(-\frac{S}{2}\theta\theta^{\top}\right)\right]$

Gaussian distribution $q(\theta) := \mathcal{N}(\theta | m, S^{-1})$ Natural parameters $\lambda := \{Sm, -S/2\}$ Expectation parameters $\mu := \{\mathbb{E}_q(\theta), \mathbb{E}_q(\theta \theta^{\top})\}$

Step 2: Use the Bayesian Learning Rule to optimize

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

Natural gradient (NatGrad)
Step 2: Use the Bayesian Learning Rule to optimize

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

Natural gradient (NatGrad)

Optimality condition:

$$\nabla_{\mu} H(q_*) = \nabla_{\mu} \mathbb{E}_{q_*} [\mathscr{E}(\theta)]$$

Step 2: Use the Bayesian Learning Rule to optimize

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

Natural gradient (NatGrad)

Optimality condition:

$$\nabla_{\mu} H(q_*) = \nabla_{\mu} \mathbb{E}_{q_*} [\mathscr{E}(\theta)]$$
$$-\lambda_* = \nabla_{\mu} \mathbb{E}_{q_*} [\mathscr{E}(\theta)]$$

For minimal Exp-Family:

Step 2: Use the Bayesian Learning Rule to optimize

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

Natural gradient (NatGrad)

Optimality condition:

For minimal Exp-Family:

$$\nabla_{\mu} H(q_*) = \nabla_{\mu} \mathbb{E}_{q_*} [\mathscr{E}(\theta)]$$
$$-\lambda_* = \nabla_{\mu} \mathbb{E}_{q_*} [\mathscr{E}(\theta)]$$

"Information matching" due to the entropy term

- 1. Natural gradients are essential, & contain higher-order information about the loss, e.g., 1st and 2nd derivatives
- 2. This info is then assigned to appropriate natural params

Step 2: Use the Bayesian Learning Rule to optimize

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

Natural gradient (NatGrad)

Optimality condition:

For minimal Exp-Family:

$$\nabla_{\mu} H(q_*) = \nabla_{\mu} \mathbb{E}_{q_*} [\mathscr{E}(\theta)]$$
$$-\lambda_* = \nabla_{\mu} \mathbb{E}_{q_*} [\mathscr{E}(\theta)]$$

"Information matching" due to the entropy term

- 1. Natural gradients are essential, & contain higher-order information about the loss, e.g., 1st and 2nd derivatives
- 2. This info is then assigned to appropriate natural params

By changing Q & approx. to natGrads, we can choose the kind of "information" and recover many learning algorithms.

Deriving Learning-Algorithms from the Bayesian Learning Rule



Deriving Learning-Algorithms from the Bayesian Learning Rule



Deriving Learning-Algorithms from the Bayesian Learning Rule



Deriving Learning-Algorithms from the Bayesian Learning Rule



Deriving Learning-Algorithms from the Bayesian Learning Rule



Deriving Learning-Algorithms from the Bayesian Learning Rule



Deriving Learning-Algorithms from the Bayesian Learning Rule



Deriving Learning-Algorithms from the Bayesian Learning Rule



Deriving Learning-Algorithms from the Bayesian Learning Rule



Deriving Learning-Algorithms from the Bayesian Learning Rule



Deriving Learning-Algorithms from the Bayesian Learning Rule



Deriving Learning-Algorithms from the Bayesian Learning Rule



See Section 1.3.1 in Khan and Rue, 2021

Gradient Descent from Bayes

Gradient descent: $\theta \leftarrow \theta - \rho \nabla_{\theta} \ell(\theta)$

Gradient descent: $\theta \leftarrow \theta - \rho \nabla_{\theta} \ell(\theta)$

Derived by choosing Gaussian with fixed covariance

Gradient descent: $\theta \leftarrow \theta - \rho \nabla_{\theta} \ell(\theta)$

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

Derived by choosing Gaussian with fixed covariance

Gradient descent: $\theta \leftarrow \theta - \rho \nabla_{\theta} \ell(\theta)$

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

Derived by choosing Gaussian with fixed covariance

Gradient descent: $\theta \leftarrow \theta - \rho \nabla_{\theta} \ell(\theta)$

Derived by choosing Gaussian with fixed covariance

Gradient descent: $\theta \leftarrow \theta - \rho \nabla_{\theta} \ell(\theta)$

Bayes Learn Rule: $m \leftarrow m - \rho \nabla_m \ell(m)$

Derived by choosing Gaussian with fixed covariance

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

Learning Algorithm	Posterior Approx.	Natural-Gradient Approx.	Sec.
Optimization Algorithms			
Gradient Descent	Gaussian (fixed cov.)	Delta method	1.3
Newton's method	Gaussian	"	1.3
$Multimodal \ optimization \ {}_{\rm (New)}$	Mixture of Gaussians	"	3.2
Deep-Learning Algorithms			
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 scal- ing, 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) $_{(New)}$	Gaussian (diagonal cov.)	Gauss-Newton Hessian approx. in Adam & no square-root scaling	4.4
Variational OGN (New)	"	Remove delta method from OGN	4.4
BayesBiNN (New)	Bernoulli	Remove delta method from STE	4.5
Approximate Bayesian Inference Algorithms			
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 $_{(New)}$	Mixture of Exp-family	None	5.4

We can compute uncertainty using a variant of Adam.

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

Bayes leads to robust solutions

Avoiding large losses



Bayes leads to robust solutions

Avoiding large losses Avoiding sharp minima



Uncertainty of Deep Nets

VOGN: A modification of Adam but match the performance on ImageNet



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

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

Uncertainty of Deep Nets

VOGN: A modification of Adam but match the performance on ImageNet



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

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

20

RMSprop/Adam from Bayes

RMSprop

BLR for Gaussian approx

$$\begin{split} s &\leftarrow (1-\rho)s + \rho [\hat{\nabla}\ell(\theta)]^2 \\ \theta &\leftarrow \theta - \alpha (\sqrt{s}+\delta)^{-1} \hat{\nabla}\ell(\theta) \end{split}$$

$$S \leftarrow (1 - \rho)S + \rho(\boldsymbol{H}_{\boldsymbol{\theta}})$$
$$m \leftarrow m - \alpha \boldsymbol{S}^{-1} \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta})$$

20

RMSprop/Adam from Bayes

RMSprop

BLR for Gaussian approx

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

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

To get RMSprop, make the following choices

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

20

RMSprop/Adam from Bayes

RMSprop

BLR for Gaussian approx

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

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

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])

Variational Online Gauss-Newton

VOGN

RMSprop

$$g \leftarrow \hat{\nabla}\ell(\theta)$$
$$s \leftarrow (1-\rho)s + \rho g^2$$
$$\theta \leftarrow \theta - \alpha(\sqrt{s} + \delta)^{-1}g$$

$$g \leftarrow \hat{\nabla}\ell(\theta), \text{ where } \theta \sim \mathcal{N}(m, \sigma^2)$$
$$s \leftarrow (1-\rho)s + \rho(\Sigma_i g_i^2)$$
$$m \leftarrow m - \alpha(s+\gamma)^{-1} \nabla_{\theta}\ell(\theta)$$
$$\sigma^2 \leftarrow (s+\gamma)^{-1}$$

import torch
+import torchsso

train_loader = torch.utils.data.DataLoader(train_dataset)
model = MLP()

```
-optimizer = torch.optim.Adam(model.parameters())
+optimizer = torchsso.optim.VOGN(model, dataset_size=len(train_loader.dataset))
```

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

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



Image Segmentation

Uncertainty (entropy of class probs)

(By Roman Bachmann)²²



Image Segmentation

Uncertainty (entropy of class probs)

(By Roman Bachmann)²²

Tuning VOGN

The trick is to mimic Adam's trajectory as closely as possible

Tuning VOGN: Currently, there is no common recipe for tuning the algorithmic hyperparameters for VI, especially for large-scale tasks like ImageNet classification. One key idea we use in our experiments is to start with Adam hyperparameters and then make sure that VOGN training closely follows an Adam-like trajectory in the beginning of training. To achieve this, we divide the tuning into an optimisation part and a regularisation part. In the optimisation part, we first tune the hyperparameters of a deterministic version of VOGN, called the online Gauss-Newton (OGN) method. This method, described in Appendix Q, is more stable than VOGN since it does not require MC sampling, and can be used as a stepping stone when moving from Adam/SGD to VOGN. After reaching a competitive performance to Adam/SGD by OGN, we move to the regularisation part, where we tune the prior precision δ , the tempering parameter τ , and the number of MC samples K for VOGN. We initialise our search by setting the prior precision δ using the L2-regularisation parameter used for OGN, as well as the dataset size N. Another technique is to warm-up the parameter τ towards $\tau = 1$ (also see the "momentum and initialisation" part). Setting τ to smaller values usually stabilises the training, and increasing it slowly also helps during tuning. We also add an *external* damping factor $\gamma > 0$ to the moving average s_t. This increases the lower bound of the eigenvalues of the diagonal covariance Σ_t and prevents the noise and the step size from becoming too large. We find that a mix of these techniques works well for the problems we considered.

Sec 3, last paragraph in Osawa et al. "Practical Deep Learning with Bayesian Principles." NeurIPS (2019). 23



Human Learning at the age of 6 months.

Comp Learning with

NEURAL INFORMATION MICCISSING STRIEMS

Bayesian Principles

by Mohammad Emtiyaz Khan · Dec 9, 2019

NeurIPS 2019 Tutorial



8.084 views · Dec 9, 2019

by <u>Weienne Sze</u> 7,163 views - Dec 9, 2019

Past and New Work

Natural Gradient Variational Inference

- 1. Khan and Lin. "Conjugate-computation variational inference: Converting variational inference in non-conjugate models to inferences in conjugate models." Alstats (2017).
- 2. Khan and Nielsen. "Fast yet simple natural-gradient descent for variational inference in complex models." (2018) ISITA.

• Mixture of Exponential family

3. Lin et al. "Fast and Simple Natural-Gradient Variational Inference with Mixture of Exponential-family Approximations," ICML (2019).

Generalization of natural gradients

- 4. Lin et al. "Handling the Positive-Definite Constraint in the Bayesian Learning Rule", ICML (2020)
- 5. Lin et al. "Tractable structured natural gradient descent using local parameterizations", ICML, (2021)
- Gaussian approx ↔ Newton-variants



Wu Lin (UBC)



Mark Schmidt (UBC)



Frank Nielsen (Sony)
Gaussian Approximation and DL

- 1. Khan, et al. "Fast and scalable Bayesian deep learning by weight-perturbation in Adam." *ICML* (2018).
- 2. Mishkin et al. "SLANG: Fast Structured Covariance Approximations for Bayesian Deep Learning with Natural Gradient" NeurIPS (2018).
- 3. Osawa et al. "Practical Deep Learning with Bayesian Principles." NeurIPS (2019).



Extensions

Binary Neural Networks (Bernoulli approx)

1. Meng, et al. "Training Binary Neural Networks using the Bayesian Learning Rule." *ICML* (2020).

Gaussian Process

2. Chang et al. "Fast Variational Learning in State-Space GP Models", MLSP (2020)

- For sparse GPs, BLR is a generalization of [1]





Roman Bachmann (Intern from EPFL)

Xiangming Meng (RIKEN-AIP)







Paul Chang (Aalto University)

W. J. Wilkinson (Aalto University) Arno Solin (Aalto University)

1. Hensman et al. "Gaussian Process for Big Data", UAI (2013)

Dual Perspective of the Bayesian Learning Rule

Memorable Examples Connections to Gaussian Process Continual learning Adaptation with K-priors

Relevance of Data Examples

Which examples are most relevant for the classifier? Red circle vs Blue circle.



Model view vs Data view

Bayes "automatically" defines data-relevance



(By Roman Bachmann) ³⁰

Model view vs Data view

Bayes "automatically" defines data-relevance



 Gaussian approx fom Bayes learning rule turn NN into Linear models & Gaussian Process (GPs) [1].

$$\sum_{i=1}^{N} \ell(y_i, f_{\theta}(x_i))$$
neural network

. .

 Gaussian approx fom Bayes learning rule turn NN into Linear models & Gaussian Process (GPs) [1].

$$\sum_{i=1}^{N} \ell(y_i, f_{\theta}(x_i)) \approx \sum_{i=1}^{N} \frac{1}{\sigma_i^2} [\tilde{y}_i - \phi_i(x_i)^{\top} \theta]^2$$

neural network

31

 Gaussian approx fom Bayes learning rule turn NN into Linear models & Gaussian Process (GPs) [1].

$$\sum_{i=1}^{N} \ell(y_i, f_{\theta}(x_i)) \approx \sum_{i=1}^{N} \frac{1}{\sigma_i^2} [\tilde{y}_i - \phi_i(x_i)^{\top} \theta]^2$$

neural network

"Dual" variables obtained from $\nabla_{\mu} \mathbb{E}_q[\ell_i(\theta)]$ (For Gaussian approx, obtained from Jacobian, residual etc.)

 Gaussian approx fom Bayes learning rule turn NN into Linear models & Gaussian Process (GPs) [1].

$$\sum_{i=1}^{N} \ell(y_i, f_{\theta}(x_i)) \approx \sum_{i=1}^{N} \frac{1}{\sigma_i^2} [\tilde{y}_i - \phi_i(x_i)^{\top} \theta]^2$$

neural network

"Dual" variables obtained from $\nabla_{\mu} \mathbb{E}_{q}[\ell_{i}(\theta)]$ (For Gaussian approx, obtained from Jacobian, residual etc.)

• σ_i^2 define the "relevance" of the data examples. We call more relevant ones the "memorable examples".

 Gaussian approx fom Bayes learning rule turn NN into Linear models & Gaussian Process (GPs) [1].

$$\sum_{i=1}^{N} \ell(y_i, f_{\theta}(x_i)) \approx \sum_{i=1}^{N} \frac{1}{\sigma_i^2} [\tilde{y}_i - \phi_i(x_i)^{\top} \theta]^2$$

neural network

"Dual" variables obtained from $\nabla_{\mu} \mathbb{E}_{q}[\ell_{i}(\theta)]$ (For Gaussian approx, obtained from Jacobian, residual etc.)

• σ_i^2 define the "relevance" of the data examples. We call more relevant ones the "memorable examples".

 Gaussian approx fom Bayes learning rule turn NN into Linear models & Gaussian Process (GPs) [1].

$$\sum_{i=1}^{N} \ell(y_i, f_{\theta}(x_i)) \approx \sum_{i=1}^{N} \frac{1}{\sigma_i^2} [\tilde{y}_i - \phi_i(x_i)^{\top} \theta]^2$$

neural network

"Dual" variables obtained from $\nabla_{\mu} \mathbb{E}_q[\ell_i(\theta)]$ (For Gaussian approx, obtained from Jacobian, residual etc.)

- σ_i^2 define the "relevance" of the data examples. We call more relevant ones the "memorable examples".
- Natural-gradients give "dual variables" (Bayes Duality)

^{1.} Khan et al. "Approximate Inference Turns Deep Networks into Gaussian Processes." NeurIPS (2019).





Continual Learning with Bayes



PingBo Pan (Intern from UT Sydney)



Siddharth Swaroop (University of Cambridge)



Runa Eschenhagen (Intern from University of Osnabruck)



Rich Turner (University of Cambridge)



Alexander Immer (Intern from EPFL)



Ehsan Abedi (Intern from EPFL)



Maciej Korzepa (Intern from DTU)

1. Khan et al. "Approximate Inference Turns Deep Networks into Gaussian Process", NeurIPS, 2019 2. Pan et al. Continual Deep Learning by Functional Regularisation of Memorable Past, NeurIPS, 2020

Continual Learning

Standard Deep Learning



Continual Learning

Standard Deep Learning



Continual Learning: past classes never revisited



Kirkpatrick, James, et al. "Overcoming catastrophic forgetting in neural networks." *Proceedings of the national academy of sciences* 114.13 (2017): 3521-3526.

Continual Learning

Standard Deep Learning



Continual Learning: past classes never revisited



Standard training leads to catastrophic forgetting.

Kirkpatrick, James, et al. "Overcoming catastrophic forgetting in neural networks." *Proceedings of the national academy of sciences* 114.13 (2017): 3521-3526.

Inaccuracy of Weight-Priors 'Add Data' task. Binary classification with base model Logistic regression (Zero offset, ie, decision boundary pass through the origin). Each task N=500, each class 250 examples.

1. Kirkpatrick et al. Overcoming catastrophic forgetting in neural networks. PNAS, 2017.

Inaccuracy of Weight-Priors



'Add Data' task.

Binary classification with Logistic regression (Zero offset, ie, decision boundary pass through the origin).

Each task N=500, each class 250 examples.

Inaccuracy of Weight-Priors



'Add Data' task.

Binary classification with Logistic regression (Zero offset, ie, decision boundary pass through the origin).

Each task N=500, each class 250 examples.

Inaccuracy of Weight-Priors



'Add Data' task.

Binary classification with Logistic regression (Zero offset, ie, decision boundary pass through the origin).

Each task N=500, each class 250 examples.





Weights Regularization [1]

$$(\boldsymbol{\theta} - \boldsymbol{\theta}_{old})^{\mathsf{T}} \Sigma_{old}^{-1} (\boldsymbol{\theta} - \boldsymbol{\theta}_{old})$$



Weights Regularization [1] $(\theta - \theta_{old})^{\mathsf{T}} \Sigma_{old}^{-1} (\theta - \theta_{old})$

Functional Regularization of memorable past (FROMP) [2]



Weights Regularization [1] $(\theta - \theta_{old})^{\top} \Sigma_{old}^{-1} (\theta - \theta_{old})$

Functional Regularization of memorable past (FROMP) [2] $KL(p(\theta)||q(\theta)) \approx KL(p(f)||q(f))$



Weights Regularization [1] $(\theta - \theta_{old})^{\top} \Sigma_{old}^{-1} (\theta - \theta_{old})$

Functional Regularization of memorable past (FROMP) [2] $KL(p(\theta)||q(\theta)) \approx KL(p(f)||q(f))$

 $\left[f(X_m) - f_{old}(X_m)\right]^\top K_{old}(X_m, X_m)^{-1} \left[f(X_m) - f_{old}(X_m)\right]$



Weights Regularization [1] $(\theta - \theta_{old})^{\top} \Sigma_{old}^{-1} (\theta - \theta_{old})$

Functional Regularization of memorable past (FROMP) [2] $KL(p(\theta)||q(\theta)) \approx KL(p(f)||q(f))$

$$\left[f(X_m) - f_{old}(X_m)\right]^\top K_{old}(X_m, X_m)^{-1} \left[f(X_m) - f_{old}(X_m)\right]$$

FROMP has the "right form" to reconstruct the "gradient of the past" — it is a type of K-prior!

Mohammad Emtiyaz Khan^{*} RIKEN Center for AI Project Tokyo, Japan emtiyaz.khan@riken.jp Siddharth Swaroop* University of Cambridge Cambridge, UK ss2163@cam.ac.uk

Abstract

Humans and animals have a natural ability to quickly adapt to their surroundings, but machine-learning models, when subjected to changes, often require a complete retraining from scratch. We present Knowledge-adaptation priors (K-priors) to reduce the cost of retraining by enabling quick and accurate adaptation for a widevariety of tasks and models. This is made possible by a combination of weight and function-space priors to reconstruct the gradients of the past, which recovers and generalizes many existing, but seemingly-unrelated, adaptation strategies. Training with simple first-order gradient methods can often recover the exact retrained model to an arbitrary accuracy by choosing a sufficiently large memory of the past data. Empirical results confirm that the adaptation can be cheap and accurate, and a promising alternative to retraining.



Joint work with Siddharth Swaroop University of Cambridge, UK







K-prior Construction

Combine weight and function-space divergences

Weight-space Function-space $\mathcal{K}(w) = \tau \mathbb{D}_w(w \| w_*) + \mathbb{D}_f(\mathbf{f}(w) \| \mathbf{f}(w_*))$

K-prior Construction

Combine weight and function-space divergences



K-prior Construction

Combine weight and function-space divergences



Faithful Gradient Reconstruction


Faithful Gradient Reconstruction



Faithful Gradient Reconstruction



No labels required, so \mathcal{M} can include any inputs!

Faithful Gradient Reconstruction



No labels required, so \mathcal{M} can include any inputs!



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.

Immer et al., Scalable Marginal Likelihood Estimation for Model Selection in Deep Learning, ICML, 2021.

• Three questions

• Three questions

- Q1: What do we know? (model)

- Three questions
 - Q1: What do we know? (model)
 - Q2: What do we not know? (uncertainty)

- Three questions
 - Q1: What do we know? (model)
 - Q2: What do we not know? (uncertainty)
 - Q3: What do we need to know? (action & exploration)

- Three questions
 - Q1: What do we know? (model)
 - Q2: What do we not know? (uncertainty)
 - Q3: What do we need to know? (action & exploration)
- Posterior approximation is the key

- Three questions
 - Q1: What do we know? (model)
 - Q2: What do we not know? (uncertainty)
 - Q3: What do we need to know? (action & exploration)
- Posterior approximation is the key
 - (Q1) Models == representation of the world

- Three questions
 - Q1: What do we know? (model)
 - Q2: What do we not know? (uncertainty)
 - Q3: What do we need to know? (action & exploration)
- Posterior approximation is the key
 - (Q1) Models == representation of the world
 - (Q2) Posterior approximations == representation of the model

- Three questions
 - Q1: What do we know? (model)
 - Q2: What do we not know? (uncertainty)
 - Q3: What do we need to know? (action & exploration)
- Posterior approximation is the key
 - (Q1) Models == representation of the world
 - (Q2) Posterior approximations == representation of the model
 - (Q3) Use posterior approximations for knowledge representation, transfer, and collection.

Approximate Bayesian Inference Team



Emtiyaz Khan Team Leader



Pierre Alquier Research Scientist



Gian Maria Marconi Postdoc



Thomas Möllenhoff Postdoc

E

Wu Lin PhD Student University of British Columbia



Dharmesh Tailor Research Assistant



Fariz Ikhwantri Part-time Student Tokyo Institute of Technology



Happy Buzaaba Part-time Student University of Tsukuba



https://team-approx-

bayes.github.io/

Evgenii Egorov Remote Collaborator Skoitsch



Siddharth Swaroop Remote Collaborator University of Cambridge



Dimitri Meunier Remote Collaborator ENSAE Paris



Peter Nickl Remote Collaborator TU Darmstadt



Erik Daxberger Remote Collaborator University of Cambridge



Alexandre Piché Remote Collaborator MILA