Learning-Algorithms from Bayesian Principles

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







The Goal of My Research

"To understand the fundamental principles of learning from data and use them to develop algorithms that can learn like living beings."



Human Learning: At the age of 6 months, learning by actively and sequentially collecting limited and correlated data.

Converged at the age of 12 months



Transfer Knowledge at the age of 14 months



Human learning \neq

Humans can learn from limited, sequential, correlated data, with a clear understanding of the world.

\neq **Deep learning**

Machines require large amount of IID data, and don't really understand the world and cannot reason about it.

My current research focuses on reducing this gap!

Learning-Algorithms from Bayesian Principles

- Practical Bayesian principles
 - To design/improve/generalize learning-algorithms.
 - By computing "posterior" distribution over unknowns.
- Generalization of many existing algorithms,
 - Classical (least-squares, Newton, HMM, Kalman.. etc).
 - Deep Learning (SGD, RMSprop, Adam).
- Helps us design new algorithms
 - Connection to Gaussian Processes.
 - Reinforcement, online, continual learning, reasoning..
- Impact: Everything with one common principle.

Learning-Algorithms by Bayesian Principles

Learning by optimization: $\theta \leftarrow \theta - \rho H^{-1} \nabla_{\theta} \ell(\theta)$

Learning by Bayes:
$$\lambda \leftarrow (1 - \rho)\lambda - \rho \nabla_{\mu} \mathbb{E}_q \left[\ell(\theta) \right]$$

 $\{V^{-1}m, V^{-1}\}$

 $\{\mathbb{E}(\theta), \mathbb{E}(\theta\theta^{\top})\}\$

Natural parameters

Expectation/moment/

mean parameters

Natural and Expectation parameters of q

e.g., Gaussian distribution

$$q(\theta) := \mathcal{N}(\theta|m, V)$$

$$\exp\left[m^{\top}V^{-1}\theta - \frac{1}{2}\theta^{\top}V^{-1}\theta\right]$$

Learning by Bayes

Learning by optimization: $\theta \leftarrow \theta - \rho H^{-1} \nabla_{\theta} \ell(\theta)$

Learning by Bayes: $\lambda \leftarrow (1-\rho)\lambda - \rho \nabla_{\mu} \mathbb{E}_q \left[\ell(\theta)\right]$

Natural and Expectation parameters of q

- Alstats 2017 ICML 2017 Classical algorithms: Least-squares, Newton's method, Kalman filters, Baum-Welch, Forward-backward, etc. Bayesian inference: EM, Laplace's method, SVI, VMP.

 - Deep learning: SGD, RMSprop, Adam.
- NeurIPS 2018 Reinforcement learning: parameter-space exploration, natural policy-search.
 - Continual learning: Elastic-weight consolidation.
 - Online learning: Exponential-weight average.

ICML 2018

ISITA 2018 ICLR 2018

- NIPS 2017 Global optimization: Natural evolutionary strategies, Gaussian homotopy, continuation method & smoothed optimization.
 - List incomplete...

$$q_{\lambda}(\theta) := \mathcal{N}(m, V) \text{ Least Squares}$$

$$\lambda \leftarrow (1 - \rho)\lambda - \rho \nabla_{\mu} \mathbb{E}_{q} \left[\ell(\theta) \right] \quad \Rightarrow \lambda_{*} = \nabla_{\mu_{*}} \mathbb{E}_{q_{*}} \left[\ell(\theta) \right]$$

$$\mathbb{E}_{q} \left[\begin{pmatrix} y - X\theta \end{pmatrix}^{\top} \begin{pmatrix} y - X\theta \end{pmatrix} + \gamma \theta^{\top} \theta \end{pmatrix} := \ell(\theta)$$

$$-\mathbb{E}_{q\lambda}[\theta]^{\top} X^{\top} y + \operatorname{trace} \left[X^{\top} X \mathbb{E}_{q\lambda}[\theta \theta^{\top}] \right]$$

$$\nabla_{\mathbb{E}_{q\lambda}}[\theta] = \left(-X^{\top} y + 0 \\ X^{\top} X + \gamma I \right) = V^{-1} m$$

$$\mathbb{E}_{q\lambda}[\theta \theta^{\top}] = \left(X^{\top} X + \gamma I \right) = V^{-1}$$

$$\mathbb{E}_{q\lambda}[\theta \theta^{\top}] = \left(X^{\top} X + \gamma I \right)^{-1} X^{\top} y$$

$$[X^{\top} X + \gamma I]^{-1} X^{\top} y$$

Neural Network

 $(X^{\top}X + \gamma I)^{-1}X^{\top}y$

$$\begin{array}{c} m \leftarrow m - \rho(S + \gamma I)^{-1}g \\ S \leftarrow (1 - \rho)S + \rho H \\ \end{array}$$
 Hessian Gradient

RMSprop

Bayes with diagonal Gaussian

$$\begin{aligned} \theta &\leftarrow \mu \\ g &\leftarrow \frac{1}{M} \sum_{i} \nabla_{\theta} \log p(\mathcal{D}_{i} | \theta) \\ s &\leftarrow (1 - \beta) s + \beta g^{2} \\ \mu &\leftarrow \mu + \alpha \ \frac{g}{\sqrt{s + \delta}} \end{aligned}$$

$$\begin{aligned} \theta &\leftarrow \mu + \epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, Ns + \lambda) \\ g &\leftarrow \frac{1}{M} \sum_{i} \nabla_{\theta} \log p(\mathcal{D}_{i} | \theta) \\ s &\leftarrow (1 - \beta)s + \beta \frac{1}{M} \sum_{i} \nabla_{\theta\theta}^{2} \log p(\mathcal{D}_{i} | \theta) \\ \mu &\leftarrow \mu + \alpha \; \frac{g + \lambda \mu / N}{s + \lambda / N} \end{aligned}$$

Learning by Bayes

Learning by Bayes works for

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

- (Minimal) ExpFam (Alstats 2017).
- Some mixtures of ExtFam (see ICML 2019).
- Kernel exponential family (Upcoming).
- The principle is to choose an appropriate sufficient statistics, which then yields an approximation of the loss.

$$\ell(\theta) \approx a^{\top} \phi(\theta)$$

$$a := \nabla_{\mu} \mathbb{E}_q[\ell(\theta)]$$

Uncertainty for Logistic-Regression



Uncertainty for Deep Learning ICML 2018





Image Segmentation

Uncertainty (entropy of class probs)

(By Roman Bachmann)¹⁵

NeurIPS 2019 Practical DL with Bayes (on ImageNet)

State-of-the-art performance and convergence rate, while preserving benefits of Bayesian principles ("well-calibrated" uncertainty).



NeurIPS 2019

17

Out-of-Distributions Test

Our method (in red) is confident on "in-distribution" data, and not overconfident on "out-of-distribution" data.



Deep Reinforcement Learning



Ruckstriesh et.al.2010, Fortunato et.al. 2017, Plapper et.al. 2017

NN Training as Inference in LinReg $(X^{T}X + \gamma I)^{-1}X^{T}y$

 $\mathbb{E}_{q} \left(\sum_{i=1}^{N} \frac{\text{likelihood}}{\operatorname{prior}} \right) + \gamma \theta^{\top} \theta \right)$ $= \sum_{i=1}^{N} \ell(y_{i}, f_{\theta}(x_{i})) + \gamma \theta^{\top} \theta$

Hessian of a linear model

$$\sum_{i=1}^{N} x_i \frac{1}{\sigma^2} x_i^T + \gamma I$$

of the last layer

Hessian of the neural-network loss



$$\begin{aligned} m \leftarrow m - \rho (S + \gamma I)^{-1} g \\ S \leftarrow (1 - \rho) S + \rho H \\ \\ \text{Hessian} \\ \end{aligned}$$
 Gradient

Gaussian approx on NN == Posterior of a linear model

$$y_i \leftarrow f_w(x_i)$$

$$\tilde{y}_i = J_i w + \epsilon_i$$



MNIST: Similarity of Examples

Kernels reveals the similarities learned by the NN. Observations are assigned higher values for correct classes.



GP Kernel

GP observations

MNIST with only 0 and 1 Digits

When trained only on 0 &1digits, we see similar patterns for 0 and 1. NN hasn't learned meaningful similarities for out-of-training classes.



GP Kernel

GP observations

MNIST with only 4 and 9 Digits

For harder tasks, the network is learns spurious correlations.



Model Selection

Marginal likelihood of the GP can be used to tune hyperparameters. In preliminary experiments, it gives better results than using ELBO!



Relevance of Examples

Given a minibatch at each iteration, we select examples with less noise (low variance of epsilon_i in the approximated linear model).



(By Roman Bachmann)

How Does This Advance AI?

Posterior Approximations are essentially representation of old data.

- eg, Gaussians represent 2nd-order statistics.

- This representation can be employed
 - To avoid forgetting (continual learning).
 - To select examples (active learning).
 - To interact with the world (reinforcement learning).
 - To intervene (causal/interpretable learning).

References

Available at https://emtiyaz.github.io/publications.html

Conjugate-Computation Variational Inference : Converting Variational Inference in Non-Conjugate Models to Inferences in Conjugate Models, (AISTATS 2017) M.E. KHAN AND W. LIN [Paper] [Code for Logistic Reg

Fast and Scalable Bayesian Deep Learning by Weight-Perturbation in Adam, (ICML 2018) M.E. KHAN, D. NIELSEN, V. TANGKARATT, W. LIN, Y. GAL, AND A. SRIVASTAVA, [ArXiv Version] [Code] [Slides]

Practical Deep Learning with Bayesian Principles, (Under review) K. Osawa, S. Swaroop, A. Jain, R. Eschenhagen, R.E. Turner, R. Yokota, M.E. Khan. [arXiv]

Approximate Inference Turns Deep Networks into Gaussian Processes, (UNDER REVIEW) M.E. KHAN, A. IMMER, E. ABEDI, M. KORZEPA. [arXiv]

A 5 page review

Fast yet Simple Natural-Gradient Descent for Variational Inference in Complex Models

Mohammad Emtiyaz Khan RIKEN Center for Advanced Intelligence Project Tokyo, Japan emtiyaz.khan@riken.jp Didrik Nielsen RIKEN Center for Advanced Intelligence Project Tokyo, Japan didrik.nielsen@riker.jp

Abstract-Bayesian inference plays an important role in advancing machine learning, but faces computational challenges when applied to complex models such as deep neural networks. Variational inference circumvents these challenges by formulating Bayesian inference as an optimization problem and solving it using gradient-based optimization. In this paper, we argue in favor of natural-gradient approaches which, unlike their gradientbased counterparts, can improve convergence by exploiting the information geometry of the solutions. We show how to derive fast yet simple natural-gradient updates by using a duality associated with exponential-family distributions. An attractive feature of these methods is that, by using natural-gradients, they are able to extract accurate local approximations for individual model components. We summarize recent results for Bayesian deep learning showing the superiority of natural-gradient approaches over their gradient counterparts.

Index Terms—Bayesian inference, variational inference, natural gradients, stochastic gradients, information geometry, exponential-family distributions, nonconjugate models. prove the rate of convergence [7]–[9]. Unfortunately, these approaches only apply to a restricted class of models known as *conditionally-conjugate* models, and do not work for non-conjugate models such as Bayesian neural networks.

This paper discusses some recent methods that generalize the use of natural gradients to such large and complex nonconjugate models. We show that, for exponential-family approximations, a duality between their natural and expectation parameter-spaces enables a simple natural gradient update. The resulting updates are equivalent to a recently proposed method called Conjugate-computation Variational Inference (CVI) [10]. An attractive feature of the method is that it naturally obtains *local* exponential-family approximations for individual model components. We discuss the application of the CVI method to Bayesian neural networks and show some recent results from a recent work [11] demonstrating





Emtiyaz Khan: Fast yet Simple Natural-Gradient Descent for Variational Inference

Acknowledgements

Slides, papers, & code are at emtiyaz.github.io



Wu Lin (Past: RA)



Zuozhu Liu

(Intern from SUTD)



Nicolas Hubacher (Past: RA)



Masashi Sugiyama Voot Tangkaratt (Director RIKEN-AIP) (Postdoc, RIKEN-AIP)

RAIDEN





Mark Schmidt (UBC)





Reza Babanezhad (UBC)



Yarin Gal (UOxford)



Akash Srivastava (UEdinburgh)

External Collaborators

Acknowledgements

Slides, papers, & code are at emtiyaz.github.io



Kazuki Osawa (Tokyo Tech)











Rich Turner (University of Cambridge)

Rio Yokota (Tokyo Tech)

Anirudh Jain (Intern from IIT-ISM, India)

Runa Eschenhagen (Intern from University of Osnabruck)





Alexander Immer (Intern from EPFL)



Ehsan Abedi (Intern from EPFL)



Maciej Korzepa (Intern from TU Denmark)



Pierre Alguier (RIKEN AIP)

Approximate Bayesian Inference Team

Looking for interns, research assistants, post-docs, and collaborators.

