# **Approximate Bayesian Inference Team Mohammad Emtiyaz Khan**



## Goals and Challenges

**Goal:** To design AI that can continually learn using Bayesian principles.

**Examples:** Uncertainty: Knowing how much we don't know, is useful to design

- Robots that can understand and reason about their environments.
- Methods that improve performance of deep-learning methods.



**Challenge:** Computation of the posterior distribution is difficult

Main Idea: Approximate integration by using optimization, and design simple

algorithms that can be implemented within existing deep learning frameworks

## Fast and Simple Algorithms for Variational Inference

Variational InferenceParameters
$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{\int p(\mathcal{D}|\theta)p(\theta)d\theta}$$
DataVariational Approximation $\approx q_{\lambda}(\theta) = \operatorname{ExpFamily}(\lambda)$ Maximize the Evidence Lower Bound (ELBO): $\max_{\lambda} \mathcal{L}(\lambda) := \mathbb{E}_{q_{\lambda}} \left[ \log p(\mathcal{D}, \theta) - \log q_{\lambda}(\theta) \right]$ 

### VI with Natural-Gradient Descent

Sato 2001, Honkela et al. 2010, Hoffman et.al. 2013

NGD: 
$$\lambda \leftarrow \lambda + \rho F(\lambda)^{-1} \nabla_{\lambda} \mathcal{L}$$
 Natural Gradient

Fisher Information Matrix (FIM)

$$F(\lambda) := \mathbb{E}_{q_{\lambda}} \left[ \nabla \log q_{\lambda}(\theta) \nabla \log q_{\lambda}(\theta)^{\top} \right]$$

- Fast convergence due to optimization in Riemannian manifold (not Euclidean space).

#### **Expectation** Parameters

Expectation/moment/ mean parameters

 $\mu := \mathbb{E}_{q_{\lambda}}[\phi(\theta)]$ Sufficient statistics

For Gaussians, it's mean and correlation matrix  $\mathbb{E}_{q_{\lambda}}[\theta] = m \qquad \mathbb{E}_{q_{\lambda}}[\theta\theta^{\top}] = mm^{\top} + V$ 

A key relationship:  $F(\lambda)^{-1} 
abla_{\lambda} \mathcal{L} = 
abla_{\mu} \mathcal{L}$ Natural Gradient wrt Gradient wrt expectation natural parameter parameter

### Gradient descent (GD) : $\lambda \leftarrow \lambda + \rho \nabla_{\lambda} \mathcal{L}$

But requires additional computations.

Can we simplify/reduce this computation?

NGD:  $\lambda \leftarrow \lambda + \rho \nabla_{\mu} \mathcal{L}$ 

**Example: Linear Regression**  $q_{\lambda}(\theta) := \mathcal{N}(m, V)$  $\mathbb{E}_{q}\left[(y - X\theta)^{\top}(y - X\theta) + \gamma\theta^{\top}\theta - \log q_{\lambda}(\theta)\right]$  $-\mathbb{E}_{q_{\lambda}}[\theta]^{\top}X^{\top}y + \operatorname{trace}\left[X^{\top}X\mathbb{E}_{q_{\lambda}}[\theta\theta^{\top}]\right]$  $\nabla_{\mathbb{E}_{q_{\lambda}}[\theta]} = \begin{pmatrix} -X^{\top}y & + & 0 & - & V^{-1}m \\ \nabla_{\mathbb{E}_{q_{\lambda}}[\theta\theta^{\top}]} = \begin{pmatrix} X^{\top}X & + & \gamma I & - & V^{-1} \end{pmatrix}$  $m \leftarrow (1 - \rho)m - \rho \left[ X^{\top} X + \gamma I \right]^{-1} X^{\top} y$ 

## **Bayesian Neural Network**

$$\mathbb{E}_{q} \left( \sum_{i=1}^{N} \frac{\text{likelihood}}{\log p(y_{i} | f_{\theta}(x_{i})) + \gamma \theta^{\top} \theta} - \log q_{\lambda}(\theta) \right)$$

$$\begin{split} m \leftarrow m - \beta (S + \gamma I)^{-1} [g_i(\theta) + \gamma m] \\ S \leftarrow (1 - \beta) S + \beta H_i(\theta) & \xrightarrow{\text{Back-propagated} \\ \text{gradient & Hessian}} \end{split}$$

 $\theta \sim q_{\lambda}(\theta), \qquad g_i(\theta) := -\nabla_{\theta} \log p(y_i | f_{\theta}(x_i)),$  $V^{-1} \leftarrow S + \gamma I, \qquad H_i(\theta) := -\nabla_\theta^2 \log p(y_i | f_\theta(x_i))$ 

MLE vs NGD-VI	
RMSprop for MLE	NGD for mean-field VI
$\begin{aligned} \theta &\leftarrow \mu \\ g &\leftarrow \frac{1}{M} \sum_{i} \nabla_{\theta} \log p(\mathcal{D}_{i}   \theta) \\ s &\leftarrow (1 - \beta) s + \beta g^{2} \\ \mu &\leftarrow \mu + \alpha \; \frac{g}{\sqrt{s + \delta}} \end{aligned}$	$ \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} \left[ \nabla_{\theta} \log p(\mathcal{D}_{i} \theta) \right]^{2} \\ \mu &\leftarrow \mu + \alpha \; \frac{g + \lambda \mu / N}{s + \lambda / N} \end{aligned} $
Variational Online Gauss-Newton (VOGN) $s \leftarrow (1 - \beta)s + \beta \frac{1}{M} \sum \nabla_{\theta \theta}^2 \log p(\mathcal{D}_i   \theta)$	

Variational RMSprop (Vprop)  $s \leftarrow (1 - \beta)s + \beta g^2$ 



#### References

- Fast and Scalable Bayesian Deep Learning by Weight-Perturbation in Adam, (ICML 2018), Khan, Nielsen, Tangkaratt, Lin, Gal, and Srivastava.
- SLANG: Fast Structured Covariance Approximations for Bayesian Deep Learning with Natural Gradient, (NeurIPS 2018), Mishkin, Kunstner, Nielsen, Schmidt, Khan.
- Fast and Simple Natural-Gradient Descent for Variational Inference in Complex Models (ISITA 2018), Khan and Nielsen. 3