

# Variational Inference based on Robust Divergences

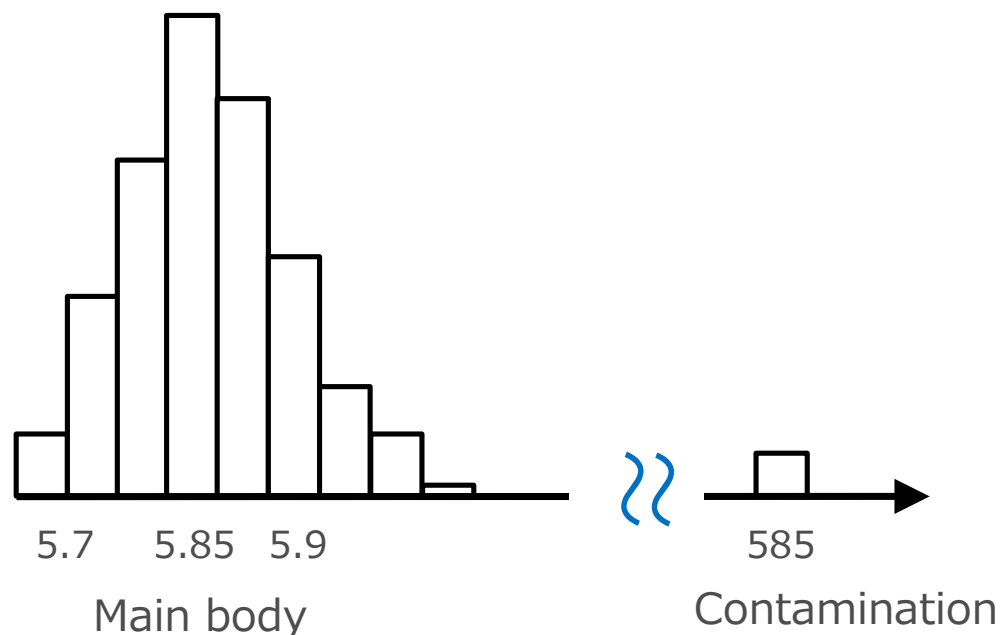
Futoshi Futami (1, 2)   Issei Sato (1, 2)   Masashi Sugiyama (2, 1)

(1) The University of Tokyo   (2) RIKEN

Advances in Approximate Bayesian Inference

December 8, 2017

# What is the outlier-robust inference?



Samples are generated from some unknown distribution.

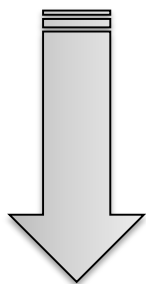
$$\{x_i\}_{i=1}^N \sim p^*(x) \quad p^*(x) = (1 - \varepsilon)p_0^*(x) + \varepsilon\delta(x)$$

Main body                      Contamination

We aim at placing an **estimated probability distribution close to the main body of the unknown distribution.**

- Estimate  $p^*(x)$  by using  $p(x; \theta)$  .
- Generalization error is measured by KL divergence:

$$D_{\text{KL}}(p^*(x) \| p(x; \theta)) = \int p^*(x) \log \left( \frac{p^*(x)}{p(x; \theta)} \right) dx.$$



Empirical approximation

$$\hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \delta(x, x_i),$$

$$\arg \min_{\theta} D_{\text{KL}}(\hat{p}(x) \| p(x; \theta)) \quad \longrightarrow \quad 0 = \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial \theta} \ln p(x_i; \theta).$$

Maximum likelihood estimation **is sensitive to outliers**  
**because it treats all data points equally.**

# Robust divergences

$\beta$  divergence (density power divergence)

$$D_{\beta}(g\|f) = \frac{1}{\beta} \int g(x)^{1+\beta} dx + \frac{\beta+1}{\beta} \int g(x)f(x)^{\beta} dx + \int f(x)^{1+\beta} dx$$

(Basu et al. [1998])

$\gamma$  divergence


$$D_{\gamma}(g\|f) = \frac{1}{\gamma(1+\gamma)} \ln \int g(x)^{1+\gamma} dx - \frac{1}{\gamma} \ln \int g(x)f(x)^{\gamma} dx + \frac{1}{1+\gamma} \ln \int f(x)^{1+\gamma} dx$$

(Fujisawa and Eguchi. [2008])

# Robust divergence minimization

Minimizing empirical  $\beta$  or  $\gamma$  divergence instead of KL

$$\arg \min_{\theta} D_{\beta} (\hat{p}(x) \| p(x; \theta))$$



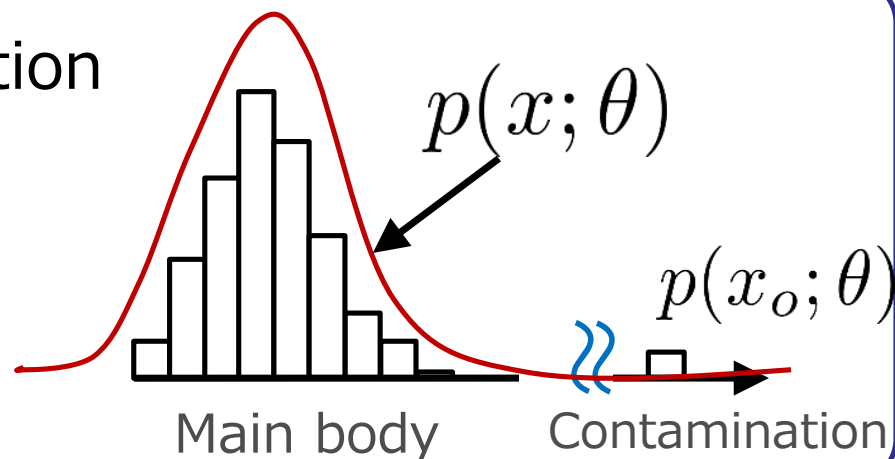
$$0 = \frac{1}{N} \sum_{i=1}^N \underline{p(x_i; \theta)^{\beta}} \frac{\partial}{\partial \theta} \ln p(x_i; \theta) - \mathbb{E}_{p(x; \theta)} \left[ p(x; \theta)^{\beta} \frac{\partial}{\partial \theta} \ln p(x; \theta) \right]$$

Density power weights

The likelihood weighted according to the power of the probability for each data point.

We want to model the distribution of the main body of data.

Outliers  $x_o$  have small  $p(x_o; \theta)$ .



# Bayesian inference (reformulation)

5

$\theta$  : random variable  
 $p(\theta)$ :prior distribution

Bayes' theorem



$$p(\theta|x_{1:N}) = \frac{p(x_{1:N}|\theta)p(\theta)}{p(x_{1:N})}$$

Reformulation of Bayesian inference

Zellner[1988], Zhu et al. [2014]

$$\arg \min_{q(\theta) \in \mathcal{P}} L(q(\theta))$$

$$L(q(\theta)) = D_{\text{KL}}(q(\theta) \| p(\theta)) - \int q(\theta) (-N d_{\text{KL}}(\hat{p}(x) \| p(x|\theta))) d\theta$$

Cross entropy:  $d_{\text{KL}}(\hat{p}(x) \| p(x|\theta)) = -\frac{1}{N} \sum_{i=1}^N \ln p(x_i|\theta)$

Solution

$$q(\theta) = \frac{e^{-N d_{\text{KL}}(\hat{p}(x) \| p(x|\theta))} p(\theta)}{\int e^{-N d_{\text{KL}}(\hat{p}(x) \| p(x|\theta))} p(\theta) d\theta}$$



$$p(\theta|x_{1:N}) = \frac{p(x_{1:N}|\theta)p(\theta)}{p(x_{1:N})}$$

# Variational inference

$$\arg \min_{q(\theta) \in \mathcal{P}} L(q(\theta))$$

$$L(q(\theta)) = D_{\text{KL}}(q(\theta) \| p(\theta)) - \int q(\theta) (-N d_{\text{KL}}(\hat{p}(x) \| p(x|\theta))) d\theta$$

This is often intractable analytically, we need some approximation method.

Restrict the domain of the optimization problem to analytically tractable distributions

$$q(\theta; m) \in \mathcal{Q}$$

$$\arg \min_{\underline{q(\theta; m) \in \mathcal{Q}}} L(q(\theta; m))$$

- This method is called **variational inference**.
- $-L(q(\theta; \lambda))$  is called the **evidence lower-bound (ELBO)**.

# Summary up to now

## Maximum Likelihood estimation

$$\arg \min_{\theta} D_{\text{KL}}(\hat{p}(x) \| p(x; \theta))$$



## Robust estimation

$$\arg \min_{\theta} D_{\beta}(\hat{p}(x) \| p(x; \theta))$$

## Bayesian Inference

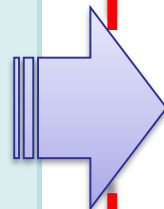
$$\arg \min_{q(\theta) \in \mathcal{P}} L(q(\theta)) = D_{\text{KL}}(q(\theta) \| p(\theta))$$



## Variational inference

$$\arg \min_{q(\theta; m) \in \mathcal{Q}} L(q(\theta)) = D_{\text{KL}}(q(\theta) \| p(\theta))$$

$$- \int q(\theta) (-N d_{\text{KL}}(\hat{p}(x) \| p(x|\theta))) d\theta$$



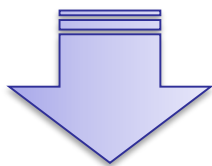
How to incorporate robust divergence property ?



# Interpretation

## Bayesian Inference

$$\arg \min_{q(\theta) \in \mathcal{P}} \underbrace{D_{\text{KL}}(q(\theta) \| p(\theta)) - \int q(\theta) (-N d_{\text{KL}}(\hat{p}(x) \| p(x|\theta))) d\theta}_{\text{seems like...}}$$



$$\arg \min_{q(\theta) \in \mathcal{P}} \mathbb{E}_{q(\theta)} [D_{\text{KL}}(\hat{p}(x) \| p(x|\theta))] + \frac{1}{N} D_{\text{KL}}(q(\theta) \| p(\theta)).$$

## Maximum Likelihood estimation

$$\arg \min_{\theta} D_{\text{KL}}(\hat{p}(x) \| p(x; \theta))$$

$$\arg \min_{q(\theta) \in \mathcal{P}}$$

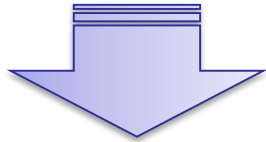
Expected  
likelihood

+

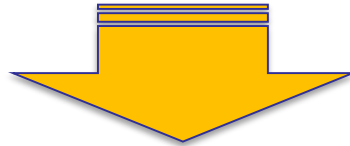
Regularization:  
close to prior

seems like...

$$\arg \min_{q(\theta) \in \mathcal{P}} \underbrace{\mathbb{E}_{q(\theta)} [D_{\text{KL}} (\hat{p}(x) \| p(x|\theta))]} + \frac{1}{N} D_{\text{KL}} (q(\theta) \| p(\theta)) .$$



$$\arg \min_{q(\theta) \in \mathcal{P}} \underbrace{\mathbb{E}_{q(\theta)} [D_{\beta} (\hat{p}(x) \| p(x|\theta))]} + \frac{1}{N} D_{\text{KL}} (q(\theta) \| p(\theta)) .$$



$$\arg \min_{q(\theta) \in \mathcal{P}} L_{\beta}(q(\theta)),$$

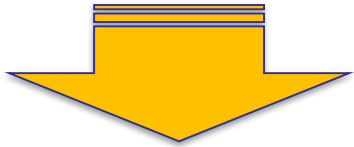
$$L_{\beta}(q(\theta)) = D_{\text{KL}}(q(\theta) \| p(\theta)) - \int q(\theta) (-N d_{\beta} (\hat{p}(x) \| p(x|\theta))) .$$

$$\beta \text{ Cross entropy: } d_{\beta}(\hat{p}(x) \| p(x|\theta)) = -\frac{\beta+1}{\beta} \frac{1}{N} \sum_{i=1}^N p(x_i|\theta)^{\beta} + \int p(x|\theta)^{1+\beta} dx .$$

$$\arg \min_{q(\theta) \in \mathcal{P}} L_{\beta}(q(\theta)),$$

$$L_{\beta}(q(\theta)) = D_{\text{KL}}(q(\theta) \| p(\theta)) - \int q(\theta) (-N d_{\beta}(\hat{p}(x) \| p(x|\theta))) d\theta.$$

$$\beta \text{ Cross entropy: } d_{\beta}(\hat{p}(x) \| p(x|\theta)) = -\frac{\beta+1}{\beta} \frac{1}{N} \sum_{i=1}^N p(x_i|\theta)^{\beta} + \int p(x|\theta)^{1+\beta} dx.$$



$$q(\theta) = \frac{e^{-N d_{\beta}(\hat{p}(x) \| p(x|\theta))} p(\theta)}{\int e^{-N d_{\beta}(\hat{p}(x) \| p(x|\theta))} p(\theta) d\theta}$$

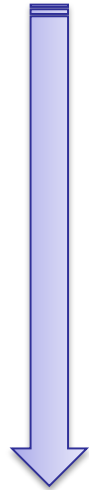
➡ Kinds of “pseudo posterior”

Conjugate relation is broken in this formulation.  
Analytical solution is intractable.

$$\arg \min_{q(\theta) \in \mathcal{P}} L_{\beta}(q(\theta)),$$

$$L_{\beta}(q(\theta)) = D_{\text{KL}}(q(\theta) \| p(\theta)) - \int q(\theta) (-N d_{\beta}(\hat{p}(x) \| p(x|\theta))) d\theta.$$

$$\beta \text{ Cross entropy: } d_{\beta}(\hat{p}(x) \| p(x|\theta)) = -\frac{\beta+1}{\beta} \frac{1}{N} \sum_{i=1}^N p(x_i|\theta)^{\beta} + \int p(x|\theta)^{1+\beta} dx.$$



Let us use **variational inference** by restricting the domain of the optimization.

$$q(\theta; m) \in \mathcal{Q}$$

**Robust variational inference**

$$\arg \min_{\underline{q(\theta; \lambda)} \in \mathcal{Q}} L_{\beta}(q(\theta; m))$$

# Proposing method

## Maximum Likelihood estimation

$$\arg \min_{\theta} D_{\text{KL}}(\hat{p}(x) \| p(x; \theta))$$

$\theta$ : random variable

## Bayesian Inference

$$\arg \min_{q(\theta) \in \mathcal{P}} L(q(\theta)) = D_{\text{KL}}(q(\theta) \| p(\theta)) - \int q(\theta) (-N d_{\text{KL}}(\hat{p}(x) \| p(x|\theta))) d\theta$$

## Variational inference

$$\arg \min_{q(\theta; \lambda) \in \mathcal{Q}} L(q(\theta)) = D_{\text{KL}}(q(\theta) \| p(\theta)) - \int q(\theta) (-N d_{\text{KL}}(\hat{p}(x) \| p(x|\theta))) d\theta$$

## Robust estimation

$$\arg \min_{\theta} D_{\beta}(\hat{p}(x) \| p(x; \theta))$$

$\theta$ : random variable

## Robust Inference

$$\arg \min_{q(\theta) \in \mathcal{P}} L_{\beta}(q(\theta)) = D_{\text{KL}}(q(\theta) \| p(\theta)) - \int q(\theta) (-N d_{\beta}(\hat{p}(x) \| p(x|\theta))) d\theta$$

## Variational inference

$$\arg \min_{q(\theta; \lambda) \in \mathcal{Q}} L_{\beta}(q(\theta)) = D_{\text{KL}}(q(\theta) \| p(\theta)) - \int q(\theta) (-N d_{\beta}(\hat{p}(x) \| p(x|\theta))) d\theta$$

- We can analyze the robustness through **IFs**.
- IFs represent **relative bias of a estimated static caused by outliers**.

Empirical distribution :  $G(x) = \frac{1}{n} \sum_{i=1}^n \delta(x, x_i)$

Contaminated version of  $G$  at  $z$  :  $G_{\varepsilon,z}(x) = (1 - \varepsilon)G(x) + \varepsilon\delta(x, z)$

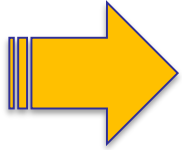
$\varepsilon$  :contamination proportion

For a static  $T$  and empirical distribution  $G$  , IF at point  $z$  is defined as:

$$\text{IF}(z, T, G) = \left. \frac{\partial}{\partial \varepsilon} T(G_{\varepsilon,z}(x)) \right|_{\varepsilon=0} = \lim_{\varepsilon \rightarrow 0} \frac{T(G_{\varepsilon,z}(x)) - T(G(x))}{\varepsilon}.$$

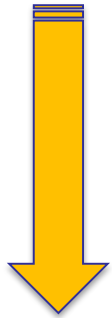
# How to use IFs ?

1. Investigate whether  $\sup_z |\text{IF}(z, m, G)| < \infty$  or not.



If it diverges, the model can be sensitive to small contamination of data.

2. How much is the predictive distribution affected by outliers ?



What we want to know is predictive distribution.

$$p(x_{\text{test}}|x_{1:N}) = \int p(\theta|x_{1:N})p(x_{\text{test}}|\theta)d\theta \approx \int q^*(\theta)p(x_{\text{test}}|\theta)d\theta.$$

$$\frac{\partial}{\partial \epsilon} \mathbb{E}_{q^*(\theta)} [p(x_{\text{test}}|\theta)] = \frac{\partial \mathbb{E}_{q^*(\theta)} [p(x_{\text{test}}|\theta)]}{\partial m} \frac{\partial m^*(G_{\epsilon, z}(x))}{\partial \epsilon}$$

# IF of variational inference

$$m^* : \text{satisfies first order condition} \quad 0 = \left. \frac{\partial}{\partial m} L \right|_{m=m^*}$$

$$q^*(\theta) := q(\theta; m^*)$$



$T$  corresponds to the variational parameter  $m$

## IF of variational inference

- For usual variational inference,

$$\frac{\partial m^*(G_{\varepsilon, z}(x))}{\partial \varepsilon} = \left( \frac{\partial^2 L}{\partial m^2} \right)^{-1} \frac{\partial}{\partial m} \mathbb{E}_{q^*(\theta)} [D_{\text{KL}}(q^*(\theta) \| p(\theta)) + N \ln p(z|\theta)]$$

- For  $\beta$ -variational inference,

$$\frac{\partial m^*(G_{\varepsilon, z}(x))}{\partial \varepsilon} = \left( \frac{\partial^2 L_{\beta}}{\partial m^2} \right)^{-1} \frac{\partial}{\partial m} \mathbb{E}_{q^*(\theta)} \left[ D_{\text{KL}}(q^*(\theta) \| p(\theta)) + N \frac{\beta + 1}{\beta} p(z|\theta)^{\beta} - \int p(x|\theta)^{1+\beta} dx \right]$$



# IF of some specific models

- Let us investigate whether  $\sup_z |\text{IF}(z, m, G)| < \infty$ .
- Consider regression and logistic regression for **Bayesian neural networks**.

Input related outlier :  $x_o \not\sim p^*(x)$   
 Output related outlier :  $y_o \not\sim p^*(y|x)$

Behavior of  $\sup_z |\text{IF}(z, W, G)|$

| Activation function | Regression           | $\beta$ - and $\gamma$ -Regression | Classification | $\beta$ - and $\gamma$ -Classification |
|---------------------|----------------------|------------------------------------|----------------|--|
| Linear              | $(x_o : U, y_o : U)$ | $(x_o : B, y_o : B)$               | $(x_o : U)$    | $(x_o : B)$                            |
| ReLU                | $(x_o : U, y_o : U)$ | $(x_o : B, y_o : B)$               | $(x_o : U)$    | $(x_o : B)$                            |
| tanh                | $(x_o : B, y_o : U)$ | $(x_o : B, y_o : B)$               | $(x_o : B)$    | $(x_o : B)$                            |

$(x_o : U, y_o : U)$  : IF is unbounded.

$(x_o : B, y_o : U)$  : IF is bounded for input related outliers, but unbounded for output related outliers.



**IF of our proposed method is always bounded.**

# Experiments on the UCI dataset

- Neural net which has two hidden layers each with 20 units and the ReLU activation function.
- We used the re-parameterization trick with 10 MC samples.
- We determine  $\beta$  or  $\gamma$  by cross-validation. (from 0.1 to 0.9 for the experiment. We found that range from 0.1 to 0.5 is enough.)
- We added outliers to training data with proportion increased from 0% to 20%.

| Dataset                            | Outliers | KL(G)      | KL(St)            | WL          | Réyni      | BB- $\alpha$      | $\beta$    | $\gamma$                 |
|------------------------------------|----------|------------|-------------------|-------------|------------|-------------------|------------|--------------------------|
| <b>concrete</b><br>N=1030<br>D=8   | 0%       | 7.46(0.34) | 7.36(0.4)         | 8.04(1.01)  | 7.16(0.39) | 7.18(0.30)        | 7.27(0.28) | <u><b>5.53(0.48)</b></u> |
|                                    | 10%      | 8.58(0.46) | 7.63(0.52)        | 10.37(1.16) | 8.04(0.43) | 7.37(0.38)        | 7.58(0.25) | <u><b>6.20(0.74)</b></u> |
|                                    | 20%      | 9.40(1.01) | 8.37(0.70)        | 11.46(0.93) | 8.63(0.52) | 7.81(0.51)        | 8.50(0.87) | <u><b>6.85(1.15)</b></u> |
| <b>powerplant</b><br>N=9568<br>D=4 | 0%       | 4.49(0.15) | 4.46(0.16)        | 4.46(0.18)  | 4.49(0.14) | 4.41(0.13)        | 4.36(0.11) | <u><b>4.28(0.14)</b></u> |
|                                    | 10%      | 4.71(0.17) | 4.59(0.15)        | 4.81(0.23)  | 4.66(0.19) | 4.56(0.17)        | 4.41(0.16) | <u><b>4.33(0.15)</b></u> |
|                                    | 20%      | 5.12(0.26) | 4.65(0.10)        | 5.04(0.25)  | 4.82(0.23) | 4.70(0.13)        | 4.52(0.15) | <u><b>4.38(0.15)</b></u> |
| <b>protein</b><br>N=45730<br>D=9   | 0%       | 5.88(0.50) | <b>4.78(0.07)</b> | 5.77(0.56)  | 4.82(0.04) | 4.81(0.04)        | 4.87(0.05) | <u><b>4.78(0.05)</b></u> |
|                                    | 10%      | 6.14(0.03) | <b>4.84(0.06)</b> | 6.14(0.028) | 4.88(0.04) | 4.86(0.04)        | 4.96(0.06) | 4.86(0.07)               |
|                                    | 20%      | 6.14(0.03) | 4.90(0.08)        | 6.14(0.031) | 4.90(0.05) | <b>4.86(0.05)</b> | 4.97(0.06) | <u><b>4.86(0.07)</b></u> |

1. We proposed an outlier-robust pseudo-Bayesian variational method by replacing the KL divergence used for data fitting to a robust divergence.
2. We analyzed our proposed method by using influence functions analytically and numerically.
3. We confirmed usefulness of our proposed method on the UCI datasets by using Bayesian neural nets.