



Sandia  
National  
Laboratories

Exceptional service in the national interest

# A Study of Bias-Variance Trade-off in Variational Inferencing Using Delta Method

SIAM UQ 2022 : Variational Inference Bridging  
Application and Theory

Niladri Das

Thomas A. Catanach

April 14, 2022

Atlanta, GA, USA



# CONTENTS

1. Variational Inference
2. Bias-Variance Trade-off
3. Delta Method Analysis on VI
4. Examples
5. Discussions



# Variational Inference

Variational inference (VI) is a method for approximating probability densities

$$\min_{\phi} [q(\theta | \phi) \parallel p(\theta | \mathcal{D})]$$

$$\begin{aligned} &= \min_{\phi} \int q(\theta | \phi) \log \frac{q(\theta | \phi)}{p(\theta | \mathcal{D})} d\theta \\ &= \boxed{\min_{\phi} \int q(\theta | \phi) \log \frac{q(\theta | \phi)}{p(\mathcal{D} | \theta) p(\theta)} d\theta} + \log p(\mathcal{D}) \end{aligned}$$

## Motivation

- In sequential data assimilation, calculating the evidence term is tractable.
- We optimize for the VI distribution parameter over sample estimates.
- Can we study the bias and variance of those sample estimates ?
- Low variance estimate can enable larger steps in the optimization problem, leading to faster convergence.



## BIAS-VARIANCE

Let  $\mathbf{X} \in \mathbb{R}^d$  be a random variable with probability distribution  $F$  with finite mean  $\mu$  and variance  $\Sigma$ .

We are interested in  $f(\mu) \in \mathbb{R}$  having access to samples from  $F$ . What can we say about the confidence in our sample estimate of  $f(\mu)$ ?

Sampling from  $F : \mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$

Estimating  $\mu : \hat{\mu} = \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j$

Estimating  $f(\mu) : f(\hat{\mu}) = f\left(\frac{1}{n} \sum_{j=1}^n \mathbf{X}_j\right)$

What can we say about the bias and variance (confidence) in the estimate  $f(\hat{\mu})$ ?



## ESTIMATION BIAS-I

Taylor series expansion of  $f(\hat{\mu})$ :

$$f(\hat{\mu}) = f(\mu) + (\hat{\mu} - \mu)^T Df(\mu) + \frac{1}{2}(\hat{\mu} - \mu)^T Hf(\mu)(\hat{\mu} - \mu) + \text{H.O.T.}$$

**Assumption 1:**

Second order approximation of  $f(\hat{\mu})$  yields:

$$\mathbb{E}(f(\hat{\mu}) - f(\mu)) = \frac{1}{2} \mathbb{E}((\hat{\mu} - \mu)^T Hf(\mu)(\hat{\mu} - \mu))$$



## ESTIMATION BIAS-II

If we assume that  $X$  is a 1-D Gaussian random variable, we analyze the exact nature of the bias. The Taylor series expansion:

$$f(\hat{\mu}) = f(\mu) + \sum_{i=1}^{\infty} \frac{1}{i!} (\hat{\mu} - \mu)^i f^i(\mu)$$

$$\mathbb{E}(f(\hat{\mu}) - f(\mu)) = \mathbb{E} \left( \sum_{i=1}^{\infty} \frac{1}{i!} (\hat{\mu} - \mu)^i f^i(\mu) \right)$$

If  $f(x) = \log(x) \rightarrow f^1(x) = \frac{1}{x} \rightarrow f^2(x) = -\frac{1}{x^2} \rightarrow f^3(x) = +2\frac{1}{x^3} \rightarrow \dots \rightarrow f^i(x) = (-1)^{i-1} \frac{(i-1)!}{x^i}$



## ESTIMATION BIAS-IV

Let us see how the co-coefficients of the term  $\frac{\Sigma^{i/2}}{n^{i/2}\mu^i}$  look like. For  $i = \{2,4,6,8,10,12\}$  the values look like

$$\left\{-\frac{1}{2}, -\frac{3}{4}, -\frac{5}{2}, -\frac{105}{8}, -\frac{189}{2}, -\frac{3465}{4}\right\}$$

For  $i = 2$  we have,

$$\mathbb{E}(f(\hat{\mu}) - f(\mu)) = -\frac{\Sigma}{2n\mu^2}$$



## ESTIMATION VARIANCE-I

### Assumption 3:

We approximate the estimate  $f(\hat{\mu})$  using the first order term only and analyze the confidence based upon this assumption.

Since  $\mathbb{E}(\hat{\mu}) = \mu$ , the expected value of the estimate  $f(\hat{\mu})$  is  $f(\mu)$  and the variance is:

$$\mathbb{E}((\hat{\mu} - \mu)^T f'(\mu) f'(\mu)^T (\hat{\mu} - \mu)) = f'(\mu)^T \frac{1}{n} \Sigma f'(\mu)$$

We can only calculate it if we know the actual  $\mu$  and the variance  $\Sigma$  of  $F$ . Since  $\hat{\mu} = \frac{1}{n} \sum_{j=1}^n X_j$ , using C.L.T. we have:

$$\sqrt{n}(\hat{\mu} - \mu) \rightarrow \mathcal{N}(0, \Sigma)$$

Using assumption 3 and C.L.T. we have:

$$\sqrt{n}(f(\hat{\mu}) - f(\mu)) \rightarrow \mathcal{N}(0, f'(\mu)^T \Sigma f'(\mu))$$

Delta Method



## ESTIMATION VARIANCE-II

### Assumption 4:

If we do not have access to actual  $\mu$  and the variance  $\Sigma$  of  $F$  we estimate the  $f'(\mu)^T \Sigma f'(\mu)$  by  $\hat{\mu}$  and  $\hat{\Sigma}$  by:

$$\hat{\Sigma} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu}) (X_i - \hat{\mu})^T$$

Combining assumption 1 and 2 we can state different levels of confidence in our estimate of  $f(\mu)$  and say that  $f(\mu)$  lies with in:

$$f(\hat{\mu}) \pm \frac{2.58}{\sqrt{n}} (f'(\hat{\mu})^T \hat{\Sigma} f'(\hat{\mu}))^{0.5}$$

with 99% confidence.



## DELTA METHOD ON VI

The KL divergence between the variational distribution and the posterior is:

$$\begin{aligned} J &= \int q(\theta | \phi) \log \frac{q(\theta | \phi)}{p(\theta)} d\theta - \int q(\theta | \phi) \log p(\mathcal{D} | \theta) d\theta + \log \int \frac{p(\mathcal{D} | \theta) p(\theta)}{q(\theta | \phi)} q(\theta | \phi) d\theta \\ &= f(\phi) - \int \left[ \frac{q(\theta(\zeta) | \phi)}{r(\theta(\zeta) | \varphi)} \right] \log p(\mathcal{D} | \theta(\zeta)) p(\zeta) d\zeta + \log \int \left[ \frac{p(\mathcal{D} | \theta(\zeta)) p(\theta(\zeta))}{r(\theta(\zeta) | \varphi)} \right] p(\zeta) d\zeta \end{aligned}$$

where  $r(\theta | \varphi)$  is an arbitrary importance distribution

- $\zeta$  is the re-parameterized random variable.
- We assume that the first term  $f(\phi)$  can be exactly calculated and do not need to be a MC estimate.
- The remaining two terms are approximated as a MC estimate. Do we need re-parameterization or should we sample directly from the importance distribution?

$$\begin{aligned}
 & - \int \left[ \frac{q(\theta(\zeta) | \phi)}{r(\theta(\zeta) | \varphi)} \right] \log p(\mathcal{D} | \theta(\zeta)) p(\zeta) d\zeta + \log \int \left[ \frac{p(\mathcal{D} | \theta(\zeta)) p(\theta(\zeta))}{r(\theta(\zeta) | \varphi)} \right] p(\zeta) d\zeta \\
 & \approx - \frac{1}{N} \sum_{i=1}^N \left[ \frac{q(\theta(\zeta_i) | \phi)}{r(\theta(\zeta_i) | \varphi)} \right] \log p(\mathcal{D} | \theta(\zeta_i)) + \log \left( \frac{1}{M} \sum_{j=1}^M \left[ \frac{p(\mathcal{D} | \theta(\zeta_j)) p(\theta(\zeta_j))}{r(\theta(\zeta_j) | \varphi)} \right] \right)
 \end{aligned}$$

We assume that the  $\phi$  (Variational parameters) and  $\varphi$  (Importance parameters) are coupled by:

$$\varphi = \phi + \eta$$

First calculate the variance and bias of the estimates:

$$\hat{x} = -\frac{1}{N} \sum_{i=1}^N \frac{q(\theta(\zeta_i) | \phi)}{r(\theta(\zeta_i) | \phi + \eta)} \log p(\mathcal{D} | \theta(\zeta_i)) \quad \text{and} \quad \hat{y} = \frac{1}{M} \sum_{j=1}^M \frac{p(\mathcal{D} | \theta(\zeta_j)) p(\theta(\zeta_j))}{r(\theta(\zeta_j) | \phi + \eta)}$$

Now,

$$\text{Var}(\hat{x}) = \frac{1}{N} \text{Var} \left( -\frac{q(\theta(\zeta) | \phi)}{r(\theta(\zeta) | \phi + \eta)} \log p(\mathcal{D} | \theta(\zeta)) \right) = \Sigma_x / N$$

We assume that we know the  $\Sigma_x$  (How to calculate this?), then

$$\sqrt{N}(\hat{\mathcal{X}} - \mathcal{X})XX \rightarrow \mathcal{N}(0, \Sigma_x)$$

Now,

$$\text{Var}(\hat{y}) = \frac{1}{M} \text{Var} \left( \frac{p(\mathcal{D} | \theta(\zeta))p(\theta(\zeta))}{r(\theta(\zeta_j) | \phi + \eta)} \right) = \Sigma_y/M$$

We assume that we know the  $\Sigma_y$  (How to calculate this?), then

$$\sqrt{M}(\hat{y} - y)XX \rightarrow d\mathcal{N}(0, \Sigma_y)$$

Use Delta-method to calculate the variance and bias of

$$\log \left( \frac{1}{M} \sum_{j=1}^M \frac{p(\mathcal{D} | \theta(\zeta_j))p(\theta(\zeta_j))}{r(\theta(\zeta_j) | \phi + \eta)} \right)$$

We have shown in the earlier step,

$$\sqrt{M}(\hat{y} - y)XX \rightarrow d\mathcal{N}(0, \Sigma_y)$$

Using Delta method,

$$\sqrt{M} \left( \log(\hat{y}) - \log(y) \right) XX \rightarrow d\mathcal{N} \left( 0, \frac{\Sigma_y}{y^2} \cdot \right)$$

Then calculate the variance and bias of

$$-\frac{1}{N} \sum_{i=1}^N \frac{q(\theta(\zeta_i) | \phi)}{r(\theta(\zeta_i) | \phi + \eta)} \log p(\mathcal{D} | \theta(\zeta_i)) \quad \text{and} \quad \log \left( \frac{1}{M} \sum_{j=1}^M \frac{p(\mathcal{D} | \theta(\zeta_j))}{r(\theta(\zeta_j) | \phi + \eta)} \right)$$

using delta method for multiple variables.

- If M and N samples are drawn independently then, (Shouldn't this be called MSE rather than Var?)

$$\text{Var}(\hat{x} + \log(\hat{y})) = \frac{\Sigma_x}{N} + \frac{\Sigma_y}{Ny^2}$$

- If M=N samples are the same samples then, using multivariate Delta method,

$$\text{Var}(\hat{x} + \log(\hat{y})) = \frac{\Sigma_x}{N} + \frac{\Sigma_y}{Ny^2} + 2 \frac{\Sigma_{xy}}{Ny}$$

- Hence,

$$\frac{\Sigma_x}{N} + \frac{\Sigma_y}{Ny^2} - 2 \frac{\sqrt{\Sigma_x \Sigma_y}}{Ny} < \text{Var}(\hat{x} + \log(\hat{y})) < \frac{\Sigma_x}{N} + \frac{\Sigma_y}{Ny^2} + 2 \frac{\sqrt{\Sigma_x \Sigma_y}}{Ny}$$