

Lecture 10

Lecturer: Andre Wibisono

Scribe: Nhi Nguyen

1 Problem: Gaussian Variational Inference

Let $\mathcal{P}_2(\mathbb{R}^d)$ denote the set of all probability distributions on \mathbb{R}^d with finite second moment, i.e., $\mathbb{E}\|x\|^2 < \infty$, and let \mathcal{G} denote the set of all non-degenerate Gaussian distributions on \mathbb{R}^d .

$$\mathcal{P}_2(\mathbb{R}^d) = \left\{ \rho : \rho \text{ is a probability distribution on } \mathbb{R}^d \text{ with } \mathbb{E}_{x \sim \rho} [\|x\|^2] < \infty \right\}$$

$$\mathcal{G} = \text{BW}(\mathbb{R}^d) = \left\{ \rho : \rho = \mathcal{N}(m, C) \text{ for some } m \in \mathbb{R}^d, C \succ 0 \in \mathbb{R}^{d \times d} \right\}$$

Given a target distribution $\nu(x) \propto e^{-f(x)} \in \mathcal{P}_2(\mathbb{R}^d)$, the goal of Gaussian variational inference is to approximate the distribution ν by a more tractable distribution $\hat{\nu} \in \mathcal{G}$ such that

$$\hat{\nu} = \arg \min_{\rho \in \mathcal{G}} \text{KL}(\rho \| \nu)$$

2 Algorithm: Gradient Flow

In order to find $\hat{\nu}$, we want to use Gradient Flow to minimize $F(\rho) = \text{KL}(\rho \| \nu)$ over $\rho \in \mathcal{G}$.

Recall from last time that we can compute the gradients of $F(\rho) = F(m, C)$ with respect to m and C as follows:

$$\frac{\partial F(m, C)}{\partial m} = \mathbb{E}_{\mathcal{N}(m, C)} [\nabla f]$$

$$\frac{\partial F(m, C)}{\partial C} = \frac{1}{2} \left(\mathbb{E}_{\mathcal{N}(m, C)} [\nabla^2 f] - C^{-1} \right)$$

where $\rho = \mathcal{N}(m, C)$ and $\nu \propto e^{-f}$.

We claim that the following ODE is good for minimizing $F(m, C)$

$$\begin{cases} \dot{m}_t &= - \mathbb{E}_{\mathcal{N}(m_t, C_t)} [\nabla f] \\ \dot{C}_t &= 2 \left(I - C_t \mathbb{E}_{\mathcal{N}(m_t, C_t)} [\nabla^2 f] \right) \end{cases} \quad (1)$$

We call this the Bures–Wasserstein Gradient Flow (BW – GF). This ODE yields a well-defined evolution of Gaussian distributions $(\rho_t = \mathcal{N}(m_t, C_t))_{t \geq 0}$, which we may optimistically believe to be a good approximation of ν .

Theorem 1. (Lambert et al. 2022) The ODE BW – GF define $\rho_t = \mathcal{N}(m_t, C_t)$, which follow the gradient flow

$$\frac{\partial \rho_t}{\partial t} = \dot{\rho}_t = -\text{grad} F(\rho_t)$$

where $F(\rho) = \text{KL}(\rho \| \nu)$ and the gradient is with respect to the Bures-Wasserstein (BW) metric

$$W_2^2(\rho_1, \rho_2) = \|m_1 - m_2\|^2 + \text{Tr} \left(C_1 + C_2 - 2 \left(C_1^{\frac{1}{2}} C_2 C_1^{\frac{1}{2}} \right)^{\frac{1}{2}} \right)$$

Note that the gradient with respect to the BW metric defines a gradient flow on a manifold. For instance, when $d = 1$, the BW metric becomes

$$W_2^2(\rho_1, \rho_2) = (m_1 - m_2)^2 + \left(\sqrt{C_1} - \sqrt{C_2} \right)^2$$

whereas the Euclidean distance between $\rho_1 = \mathcal{N}(m_1, C_1)$ and $\rho_2 = \mathcal{N}(m_2, C_2)$ is

$$\text{Euc}(\rho_1, \rho_2) = (m_1 - m_2)^2 + (C_1 - C_2)^2$$

3 Convergence results

Now, we are interested in the convergence rate for BW – GF, i.e. the gradient flow for $F(\rho) = \text{KL}(\rho \| \nu)$ on (\mathcal{G}, W_2)

Theorem 2. (Lambert et al. 2022) Assume that $\nu(x) \propto e^{-f(x)}$ for some $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and assume that f is α -strongly convex for some $\alpha \geq 0$. In other words, $\nabla^2 f(x) \succeq \alpha I$ for all $x \in \mathbb{R}^d$. This also means $\nu \propto e^{-f}$ is α -strongly log-concave (α – SLC). Also recall that $\hat{\nu} = \arg \min_{\rho \in \mathcal{G}} \text{KL}(\rho \| \nu)$. Then, along the BW – GF that defines $\rho_t = \mathcal{N}(m_t, C_t)$:

(1) For all $t \geq 0$, we have

$$W_2^2(\rho_t, \hat{\nu}) \leq e^{-2\alpha t} W_2^2(\rho_0, \hat{\nu}).$$

That means ρ_t converges to the minimizer exponentially fast when $\alpha > 0$.

(2) For all $t \geq 0$, we have

$$\underbrace{\text{KL}(\rho_t \| \nu)}_{F(\rho_t)} - \underbrace{\text{KL}(\hat{\nu} \| \nu)}_{\min_{\rho \in \mathcal{G}} F(\rho)} \leq e^{-2\alpha t} \left(\underbrace{\text{KL}(\rho_0 \| \nu)}_{F(\rho_0)} - \underbrace{\text{KL}(\hat{\nu} \| \nu)}_{\min_{\rho \in \mathcal{G}} F(\rho)} \right)$$

That means $F(\rho_t)$ converges to the minima exponentially fast when $\alpha > 0$.

(3) If $\alpha = 0$, then for all $t \geq 0$, we have

$$\text{KL}(\rho_t \| \nu) - \text{KL}(\hat{\nu} \| \nu) \leq \frac{1}{2t} W_2^2(\rho_0, \hat{\nu})$$

4 Approximation results

Exponential contraction is great! But we are also interested in how accurate the approximation $\hat{\nu} \in \mathcal{G}$ could be to the target distribution ν .

Recall that for a posterior distribution

$$\begin{aligned}\nu(x) &= p_n(X \mid y_1, \dots, y_n) \\ &\propto p_0(x) \prod_{i=1}^n p(y_i \mid x) \quad (\text{as } Y_i \mid X = x \text{ are i.i.d.}) \\ &\propto \exp \left(-f_0(x) - \sum_{i=1}^n \ell(x, y_i) \right)\end{aligned}$$

where $p_0 \propto e^{-f_0(x)}$ is the prior distribution and $\ell(x, y) = -\log p(y \mid x)$.

That means

$$v = p_n \propto e^{-f_n}$$

where

$$f_n(x) = f_0(x) + \sum_{i=1}^n \ell(x, y_i)$$

This implies

$$\lim_{n \rightarrow \infty} \frac{1}{n} f_n(x) = \lim_{n \rightarrow \infty} \frac{1}{n} f_0(x) + \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \ell(x, y_i) = 0 + \mathbb{E}_Y [\ell(x, Y)] = g(x)$$

where $g(x) = \mathbb{E}_Y [\ell(x, Y)]$. then

$$p_n(x) \propto e^{-f_n(x)} \propto \exp \left(-n \left(\frac{1}{n} f(x) \right) \right) \stackrel{n \gg 1}{\approx} \exp(-n \cdot g(x))$$

So if we assume that the target distribution is $\nu(x) \propto \exp(-n \cdot g(x))$, and if $n \gg 1$, then under certain assumptions (HW), we would obtain

$$\text{Var}_\nu(X) = O\left(\frac{1}{n}\right),$$

i.e. ν_n looks roughly unimodal. We would then expect that Gaussian variational inference would yields a good approximation $\hat{\nu} \in \mathcal{G}$ for ν_n .

Theorem 3. (Katsevich and Rigollet 2023) Assume a target distribution $\nu(x) \propto \exp(-n \cdot g(x))$, $g : \mathbb{R}^d \rightarrow \mathbb{R}$ such that

- g has a unique global minimizer m^*

- $n \gtrsim d^3$, where d is the number of dimensions and n is the number of observations

Let $m_\nu = m_{\text{EP}} = \mathbb{E}_\nu[x]$ and $C_\nu = C_{\text{EP}} = \text{Cov}_\nu(x)$, and let m_{VI} and C_{VI} be the mean and covariance of the variational Gaussian approximation to ν . Then

$$\sqrt{n} \|m_{\text{VI}} - m_\nu\|_{\ell_2} \lesssim \left(\frac{d^3}{n}\right)^{\frac{3}{2}}$$

$$n \|C_{\text{VI}} - C_\nu\|_{\text{op}} \lesssim \frac{d^3}{n}$$

Notice that if m_{Lap} and C_{Lap} are the mean and covariance of the Laplace approximation to ν , then

$$\sqrt{n} \|m_{\text{Lap}} - m_\nu\|_{\ell_2} \lesssim \left(\frac{d^3}{n}\right)^{\frac{1}{2}}$$

$$n \|C_{\text{Lap}} - C_\nu\|_{\text{op}} \lesssim \frac{\text{Const}(d)}{n}$$

The following are plots from the paper by Lambert et al. 2022 demonstrating the effectiveness of their BW – GF algorithm on a Bayesian logistic regression problem.

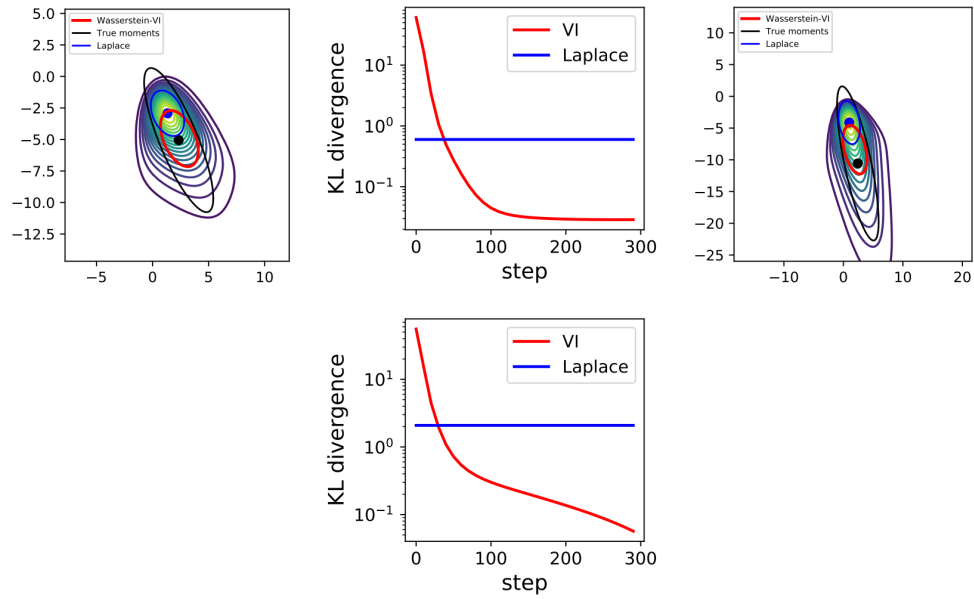


Figure 7: Results in dimension $d = 2$, $N = 10$ for a separation factor $s = 1.5$ (upper row) and $s = 2$ (lower row). The left column shows the true density via contour lines, the true mean (black dot) and covariance (black ellipsoid), and the results of the Laplace and Wasserstein VI approximations as blue and red ellipsoids respectively. The right column shows the evolution of the left KL divergence for Gaussian VI on a logarithmic scale. The corresponding KL divergence obtained with Laplace approximation is shown as a blue straight line.

References

- Katsevich, Anya and Philippe Rigollet (2023). “On the Approximation Accuracy of Gaussian Variational Inference”. In: *arXiv preprint arXiv:2301.02168*.
- Lambert, Marc et al. (2022). “Variational inference via Wasserstein gradient flows”. In: *Advances in Neural Information Processing Systems*. Ed. by Alice H. Oh et al. URL: <https://openreview.net/forum?id=K2PTuvVTF1L>.