

1 Introduction

Natural gradient method for variational inference can lead to fast convergent algorithms, but its applications are usually restricted to exponential-family approximations.

- ▶ We present a new approach to obtain natural-gradient updates for several types of approximations **outside** the class of exponential-family distributions.
- ▶ Our approach enables the derivation of **simple** updates by introducing a new type of **expectation-parameterization**.
- ▶ Our results demonstrate **faster** convergence compared to existing block-box gradient methods.

2 VI using Exp-Family Components

Given data \mathcal{D} and model $p(\mathcal{D}|\mathbf{z})$ with latent vector \mathbf{z} and prior $p(\mathbf{z})$, our goal is to approximate posterior $p(\mathbf{z}|\mathcal{D})$.

Variational inference (VI) approximates the posterior by optimizing the evidence lower bound (ELBO) \mathcal{L} induced by a variational distribution q .

Structured Approximation: We consider $q(\mathbf{w}, \mathbf{z}) = q(\mathbf{w}|\lambda_w)q(\mathbf{z}|\lambda_z)$, where

Conditional Exp-Family: $q(\mathbf{z}|\mathbf{w}, \lambda_z) := h_z(\mathbf{z}, \mathbf{w}) \exp[\langle \phi_z(\mathbf{z}, \mathbf{w}), \lambda_z \rangle - A_z(\lambda_z, \mathbf{w})]$,

Exp-Family: $q(\mathbf{w}|\lambda_w) := h_w(\mathbf{w}) \exp[\langle \phi_w(\mathbf{w}), \lambda_w \rangle - A_w(\lambda_w)]$.

We further consider the following multi-linear exponential family with N blocks.

Multi-linear Exp-Family: $q(\mathbf{z}|\lambda_1, \dots, \lambda_N) = h_z(\mathbf{z}) \exp[f(\mathbf{z}, \lambda_1, \dots, \lambda_N) - A_z(\lambda_1, \dots, \lambda_N)]$, where $f(\mathbf{z}, \lambda_1, \dots, \lambda_N)$ is a linear function w.r.t. each block λ_j given others.

Black-Box VI and Natural-Gradient VI:

$$\text{BBVI: } \lambda_z \leftarrow \lambda_z + \alpha \nabla_{\lambda_z} \mathcal{L}(\lambda_z), \quad \text{NGVI: } \lambda_z \leftarrow \lambda_z + \beta [\mathbf{F}_z(\lambda_z)]^{-1} \nabla_{\lambda_z} \mathcal{L}(\lambda_z),$$

Advantages of NGVI:

- ▶ NGVI admits a simple update in the exponential family (Khan and Lin, 2017).
- ▶ NGVI for Exp-Family: $\lambda_z \leftarrow \lambda_z + \beta \nabla_{\mathbf{m}_z} \mathcal{L}(\lambda_z)$, where \mathbf{m}_z is the expectation parameter.
- ▶ NGVI often results in faster convergence than BBVI.

Challenges of NGVI: In general, NGVI could be complicated due to the inverse of the Fisher information matrix $[\mathbf{F}_z(\lambda_z)]^{-1}$. Usually, NGVI does not admit a simple update outside the class of exponential family.

3 Simple Natural-gradient VI Update

NGVI can have a simple update in the following cases.

For a **mixture of exponential family** distribution $q(\mathbf{w}, \mathbf{z}|\lambda_w, \lambda_z)$, we define the following

- ▶ **Expectation parameters:** $\mathbf{m}_w := \mathbb{E}_{q(\mathbf{w})}[\phi_w(\mathbf{w})]$, $\mathbf{m}_z := \mathbb{E}_{q(\mathbf{w}, \mathbf{z})}[\phi_z(\mathbf{z}, \mathbf{w})]$
- ▶ **Natural parameters:** λ_w, λ_z
- ▶ **Fisher information matrix:** $\mathbf{F}_{wz}(\lambda_w, \lambda_z) = -\mathbb{E}_{q(\mathbf{w}, \mathbf{z})}[\nabla^2 \log q(\mathbf{w}, \mathbf{z}|\lambda_w, \lambda_z)]$

The following **natural gradient** update in **natural parameters:**

$$\begin{bmatrix} \lambda_w^{t+1} \\ \lambda_z^{t+1} \end{bmatrix} = \begin{bmatrix} \lambda_w^t \\ \lambda_z^t \end{bmatrix} + \beta \underbrace{\mathbf{F}_{wz}(\lambda_w^t, \lambda_z^t)^{-1}}_{\text{Natural gradient}} \begin{bmatrix} \nabla_{\lambda_w} \mathcal{L}^t \\ \nabla_{\lambda_z} \mathcal{L}^t \end{bmatrix}$$

is equivalent to

$$\text{NGVI: } \begin{aligned} \lambda_w^{t+1} &= \lambda_w^t + \beta \nabla_{\mathbf{m}_w} \mathcal{L}^t \\ \lambda_z^{t+1} &= \lambda_z^t + \beta \nabla_{\mathbf{m}_z} \mathcal{L}^t \end{aligned}$$

Similarly, for a **multi-linear exponential family** distribution, we propose to optimize λ_j given λ_{-j} . The distribution then can be re-expressed as

$$q(\mathbf{z}|\lambda_j, \lambda_{-j}) = h_z(\mathbf{z}) \exp \left[\frac{\langle \phi_j(\mathbf{z}, \lambda_{-j}), \lambda_j \rangle + f_j(\mathbf{z}, \lambda_{-j}) - A_z(\lambda_j, \lambda_{-j})}{f(\mathbf{z}, \lambda_j, \lambda_{-j})} \right]$$

For the j -th block, we define the following

- ▶ **Expectation parameters:** $\mathbf{m}_j := \mathbb{E}_{q(\mathbf{z})}[\phi_j(\mathbf{z}, \lambda_{-j})]$
- ▶ **Natural parameters:** λ_j
- ▶ **Fisher information matrix:** $\mathbf{F}_j(\lambda_j, \lambda_{-j}) = -\mathbb{E}_{q(\mathbf{z})}[\nabla_{\lambda_j}^2 \log q(\mathbf{z}|\lambda_j, \lambda_{-j})]$

The following **block natural gradient** update in **natural parameters** at block j :

$$\lambda_j^{t+1} = \lambda_j^t + \beta \underbrace{\mathbf{F}_j(\lambda_j^t, \lambda_{-j}^t)^{-1}}_{\text{Natural gradient}} \nabla_{\lambda_j} \mathcal{L}^t$$

is equivalent to

$$\text{BNGVI: } \lambda_j^{t+1} = \lambda_j^t + \beta \nabla_{\mathbf{m}_j} \mathcal{L}^t$$

The Jacobi Variant:

Instead of updating one block at each iteration, we can update all the blocks at once.

$$\text{BNGVI-J: } \lambda_j^{t+1} = \lambda_j^t + \beta \nabla_{\mathbf{m}_j} \mathcal{L}^t \quad \text{for all } j$$

This is equivalent to the following **approximate natural-gradient** update in **natural parameters:**

$$\begin{bmatrix} \lambda_1^{t+1} \\ \vdots \\ \lambda_N^{t+1} \end{bmatrix} = \begin{bmatrix} \lambda_1^t \\ \vdots \\ \lambda_N^t \end{bmatrix} + \beta \underbrace{\text{block-diag}(\mathbf{F}(\lambda_1^t, \dots, \lambda_N^t))^{-1}}_{\text{Approximate natural gradient}} \begin{bmatrix} \nabla_{\lambda_1} \mathcal{L}^t \\ \vdots \\ \nabla_{\lambda_N} \mathcal{L}^t \end{bmatrix},$$

where $\mathbf{F}(\lambda_1, \dots, \lambda_N) = -\mathbb{E}_{q(\mathbf{z})}[\nabla^2 \log q(\mathbf{z}|\lambda_1, \dots, \lambda_N)]$.

References:

- ▶ Khan and Lin. Conjugate-computation variational inference. *AISTATS*, 2017.
- ▶ Gupta et al. Shampoo: Preconditioned Stochastic Tensor Optimization *ICML*, 2018.
- ▶ Zhang et al. Noisy natural gradient as variational inference. *ICML*, 2018.

4 Examples

Example of Mixture of Exponential Family

We consider a model with a Student's t prior expressed as a scale mixture of Gaussians.

$$p(\mathcal{D}, \mathbf{z}, \mathbf{w}) = \text{InvGam}(w|a_0, a_0) \mathcal{N}(\mathbf{z}|\mathbf{0}, w\mathbf{I}) \prod_n p(\mathcal{D}_n|\mathbf{z})$$

We use the following mixture of exponential family distributions.

$$q(\mathbf{z}, \mathbf{w}) = \text{InvGam}(w|a, a) \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma}), \quad \text{where } \mathbf{z} \in \mathcal{R}^d, w \in \mathcal{R}_{++}$$

The natural parameter and expectation parameter are

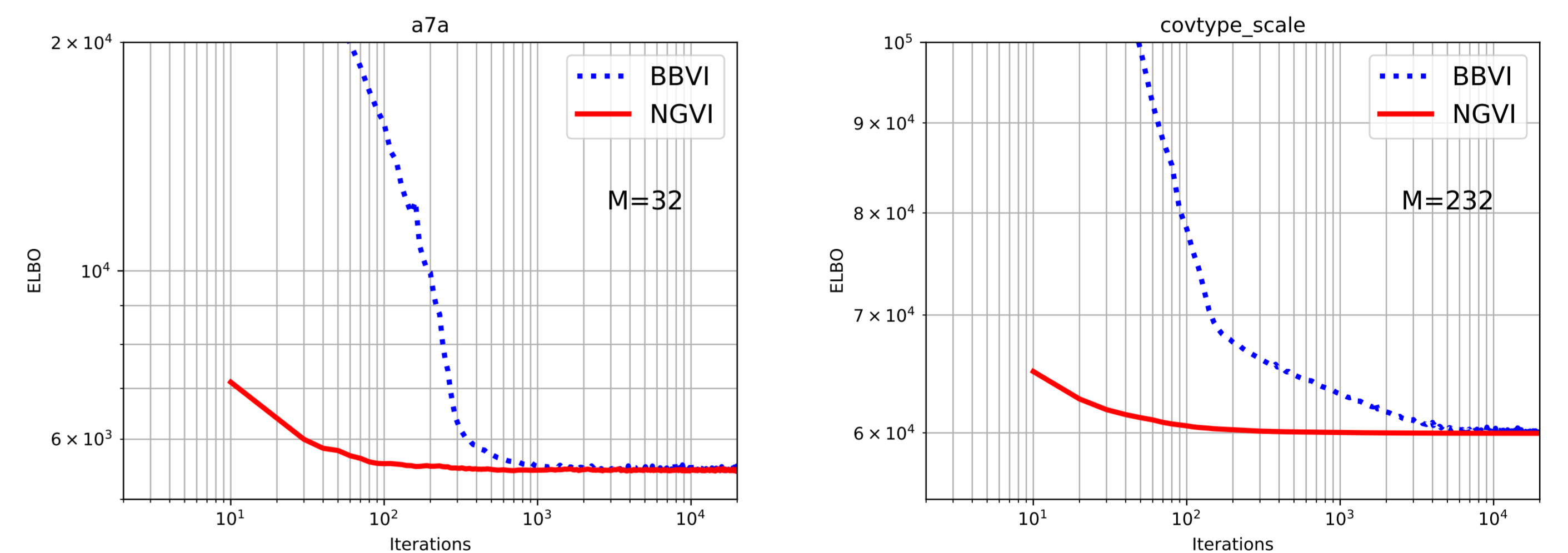
$$\begin{aligned} \lambda_z &= \left\{ \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}, -\frac{1}{2} \boldsymbol{\Sigma}^{-1} \right\}, & \mathbf{m}_z &= \{ \boldsymbol{\mu}, \boldsymbol{\mu} \boldsymbol{\mu}^T + \boldsymbol{\Sigma} \} \\ \lambda_w &= a, & \mathbf{m}_w &= -1 - (\log a - \psi(a)) \end{aligned}$$

The ELBO is $\mathcal{L} = \mathbb{E}_{q(\mathbf{z}, \mathbf{w})} \left[\sum_n \log p(\mathcal{D}_n|\mathbf{z}) + \log \frac{\mathcal{N}(\mathbf{z}|\mathbf{0}, w\mathbf{I})}{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})} + \log \frac{\text{InvGam}(w|a_0, a_0)}{\text{InvGam}(w|a, a)} \right]$.

We can re-express the update in terms of $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}^{-1}$, and a .

$$\begin{aligned} \text{NGVI: } a^{t+1} &= (1 - \beta) a^t + \beta \left(a_0 + \frac{\sum_n \nabla_a \mathbb{E}_{q(\mathbf{z}, \mathbf{w})} [\log p(\mathcal{D}_n|\mathbf{z})]}{\nabla_a \mathbb{E}_{q(\mathbf{w})} [\phi_w(\mathbf{w})]} \right) \\ (\boldsymbol{\Sigma}^{t+1})^{-1} &= (1 - \beta) (\boldsymbol{\Sigma}^t)^{-1} - 2\beta \nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{q(\mathbf{z}, \mathbf{w})} \left[\sum_n \log p(\mathcal{D}_n|\mathbf{z}) \right] + \beta \mathbf{I} \\ \boldsymbol{\mu}^{t+1} &= \boldsymbol{\mu}^t + \beta \boldsymbol{\Sigma}^{t+1} \left(\nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(\mathbf{z}, \mathbf{w})} \left[\sum_n \log p(\mathcal{D}_n|\mathbf{z}) \right] - \boldsymbol{\mu}^t \right) \end{aligned}$$

Bayesian logistic regression (t prior) with NGVI and BBVI.



Example of Mixture of Exponential Family:

We consider a model with a Gaussian prior.

$$p(\mathcal{D}, \mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{S}_0) \prod_n p(\mathcal{D}_n|\mathbf{z})$$

We use the following distribution, which is the marginal distribution of a K -mixture of Gaussians shown below. Let $\pi_K = 1 - \sum_{c=1}^{K-1} \pi_c$.

$$q(\mathbf{z}) = \sum_{w=1}^K \text{Cate}_K(w|\pi) \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w), \quad \text{where } \text{Cate}_K(w|\pi) = \exp \left(\sum_{c=1}^{K-1} \mathbb{I}_c(w) \log \frac{\pi_c}{\pi_K} + \log \pi_K \right)$$

The natural parameter and expectation parameter are

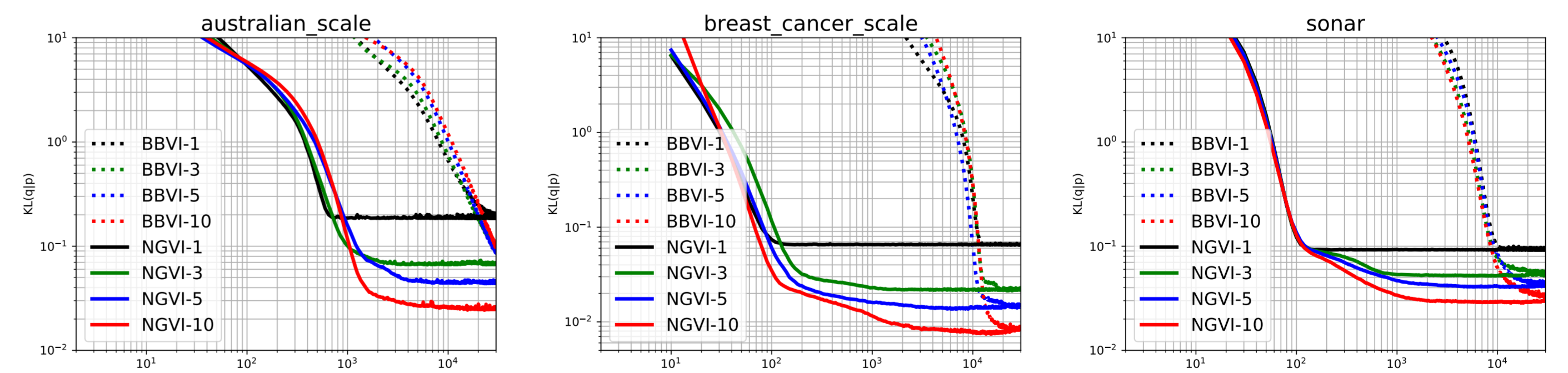
$$\begin{aligned} \lambda_z &= \left\{ \boldsymbol{\Sigma}_c^{-1} \boldsymbol{\mu}_c, -\frac{1}{2} \boldsymbol{\Sigma}_c^{-1} \right\}_{c=1}^K, & \mathbf{m}_z &= \left\{ \pi_c \boldsymbol{\mu}_c, \pi_c (\boldsymbol{\mu}_c \boldsymbol{\mu}_c^T + \boldsymbol{\Sigma}_c) \right\}_{c=1}^K \\ \lambda_w &= \left\{ \log \frac{\pi_c}{\pi_K} \right\}_{c=1}^{K-1}, & \mathbf{m}_w &= \{ \pi_c \}_{c=1}^{K-1} \end{aligned}$$

The ELBO is $\mathcal{L} = \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{z}) + \sum_n \log p(\mathcal{D}_n|\mathbf{z}) - \log q(\mathbf{z})]$.

We can re-express the update in terms of $\{\boldsymbol{\mu}_c\}_{c=1}^K$, $\{\boldsymbol{\Sigma}_c\}_{c=1}^K$, and $\{\pi_c\}_{c=1}^K$.

$$\begin{aligned} \text{NGVI: } \log \frac{\pi_c^{t+1}}{\pi_K^{t+1}} &= \log \frac{\pi_c^t}{\pi_K^t} + \beta \nabla_{\pi_c} \mathcal{L}^t \quad \text{for } c = 1, \dots, K-1 \\ (\boldsymbol{\Sigma}_c^{t+1})^{-1} &= (\boldsymbol{\Sigma}_c^t)^{-1} - \frac{2\beta}{\pi_c^t} \nabla_{\boldsymbol{\Sigma}_c} \mathcal{L}^t \quad \text{for } c = 1, \dots, K \\ \boldsymbol{\mu}_c^{t+1} &= \boldsymbol{\mu}_c^t + \frac{\beta}{\pi_c^t} \boldsymbol{\Sigma}_c^{t+1} \nabla_{\boldsymbol{\mu}_c} \mathcal{L}^t \quad \text{for } c = 1, \dots, K \end{aligned}$$

Bayesian logistic regression (Gaussian prior) with NGVI and BBVI.



Example of Multi-linear Exponential Family Approximation:

We consider a Bayesian model $p(\mathcal{D}, \mathbf{Z})$. We use the following multi-linear exponential family distribution $\mathbf{Z} \in \mathcal{R}^{d \times p}$.

$$q(\mathbf{Z}) = \mathcal{M}\mathcal{N}(\mathbf{Z}|\mathbf{W}, \mathbf{U}, \mathbf{V}), \quad \text{where } f(\mathbf{Z}, \mathbf{W}, \mathbf{U}^{-1}, \mathbf{V}^{-1}) = \text{Tr} \left(\mathbf{V}^{-1} \left(-\frac{1}{2} \mathbf{Z} + \mathbf{W} \right)^T \mathbf{U}^{-1} \mathbf{Z} \right)$$

$$\begin{aligned} \lambda_1 &= \mathbf{W}, & \lambda_2 &= \mathbf{U}^{-1}, & \lambda_3 &= \mathbf{V}^{-1} \\ \mathbf{m}_1 &= \mathbf{U}^{-1} \mathbf{W} \mathbf{V}^{-1}, & \mathbf{m}_2 &= \frac{1}{2} (\mathbf{W} \mathbf{V}^{-1} \mathbf{W}^T - \rho \mathbf{U}), & \mathbf{m}_3 &= \frac{1}{2} (\mathbf{W}^T \mathbf{U}^{-1} \mathbf{W} - d \mathbf{V}) \end{aligned}$$

Using the Gauss-Newton approximation to the Hessian matrix, we obtain the update:

$$\begin{aligned} \text{BNGVI-J: } \mathbf{W}^{t+1} &= \mathbf{W}^t + \beta \mathbf{U}^t \mathbb{E}_{q(\mathbf{Z})} [\nabla_{\mathbf{Z}} \log p(\mathcal{D}, \mathbf{Z})] \mathbf{V}^t \\ (\mathbf{U}^{t+1})^{-1} &= (\mathbf{U}^t)^{-1} + \beta \mathbb{E}_{q(\mathbf{Z})} \left[\nabla_{\mathbf{Z}} \log p(\mathcal{D}, \mathbf{Z}) \mathbf{V}^t \nabla_{\mathbf{Z}} \log p(\mathcal{D}, \mathbf{Z})^T \right] \\ (\mathbf{V}^{t+1})^{-1} &= (\mathbf{V}^t)^{-1} + \beta \mathbb{E}_{q(\mathbf{Z})} \left[\nabla_{\mathbf{Z}} \log p(\mathcal{D}, \mathbf{Z})^T \mathbf{U}^t \nabla_{\mathbf{Z}} \log p(\mathcal{D}, \mathbf{Z}) \right] \end{aligned}$$

The update is similar to Shampoo (Gupta et al., 2018). If prior $p(\mathbf{Z})$ is also a matrix-variate Gaussian distribution, the update resembles noisy K-FAC (Zhang et al. 2018).