

1 Introduction

Natural gradient method for variational inference can lead to fast convergent algorithms, but its application is usually restricted to exponential-family approximations. We extend the application to a class of distributions, and show fast faster convergence than existing block-box gradient methods.

2 VI using Exponential Family

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 $q(\mathbf{z}|\lambda_z)$.

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: $\lambda_z \leftarrow \lambda_z + \beta \nabla_{m_z} \mathcal{L}(\lambda_z)$, where \mathbf{m}_z is the expectation parameter.
- ▶ NGVI often results in faster convergence than BBVI.

Challenges of NGVI:

- ▶ NGVI could be complicated due to computing $[\mathbf{F}_z(\lambda_z)]^{-1} \nabla_{\lambda_z} \mathcal{L}(\lambda_z)$.
- ▶ $\mathbf{F}_z(\lambda_z)$ can be singular.
- ▶ Usually, NGVI does not admit a simple update outside the class of exp-family.

3 Simple Natural-gradient VI Update

Many existing works assume q to be exponential family (e.g., Gaussian). In this work, we consider more flexible approximations that are a type of mixture approximations. We consider a new NGVI update which admits a simple update in the following cases.

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

$$\text{Conditional Exp-family: } q(\mathbf{z}|\mathbf{w}, \lambda_z) := h_z(\mathbf{z}, \mathbf{w}) \exp[\langle \phi_z(\mathbf{z}, \mathbf{w}), \lambda_z \rangle - A_z(\lambda_z, \mathbf{w})], \\ q(\mathbf{w}|\lambda_w) := h_w(\mathbf{w}) \exp[\langle \phi_w(\mathbf{w}), \lambda_w \rangle - A_w(\lambda_w)].$$

We assume the set of parameters for $q(\mathbf{w})$ and $q(\mathbf{z}|\mathbf{w})$ denoted by Ω_w and Ω_z are open respectively.

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

- ▶ **Expectation** parameter: $\mathbf{m}_w := \mathbb{E}_{q(\mathbf{w})}[\phi_w(\mathbf{w})] \in \mathcal{M}_w$, $\mathbf{m}_z := \mathbb{E}_{q(\mathbf{w}, \mathbf{z})}[\phi_z(\mathbf{z}, \mathbf{w})] \in \mathcal{M}_z$
- ▶ **Natural** parameter: $\lambda = (\lambda_w, \lambda_z) \in \Omega_w \times \Omega_z$
- ▶ **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**:

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

When $\mathbf{F}_{wz}(\lambda_w^t, \lambda_z^t)$ is invertible, the update is equivalent to

$$\text{NGVI: } \lambda_w^{t+1} = \lambda_w^t + \beta \nabla_{m_w} \mathcal{L}^t \\ \lambda_z^{t+1} = \lambda_z^t + \beta \nabla_{m_z} \mathcal{L}^t$$

Now, we give a sufficient condition when $\mathbf{F}_{wz}(\lambda_w^t, \lambda_z^t)$ is invertible.

Definition 1: Minimal Conditional Exp-family (MCEF)

A conditional exp-family is said to have a minimal representation when $\mathbf{m}_w(\cdot) : \Omega_w \rightarrow \mathcal{M}_w$ and $\mathbf{m}_z(\cdot, \lambda_w) : \Omega_z \rightarrow \mathcal{M}_z$ are both one-to-one, $\forall \lambda_w \in \Omega_w$.

Theorem 1

For an MCEF representation, the FIM $\mathbf{F}_{wz}(\lambda)$ is positive-definite and invertible for all $\lambda \in \Omega$.

We further generalize 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 we assume $f(\mathbf{z}, \lambda_1, \dots, \lambda_N)$ is a linear function w.r.t. each block λ_j given others, and the set of parameters for $q(\mathbf{z})$ denoted by Ω_j is open for each block.

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

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

For the j -th block, we define the following

- ▶ **Expectation** parameter: $\mathbf{m}_j := \mathbb{E}_{q(\mathbf{z})}[\phi_j(\mathbf{z}, \lambda_{-j})] \in \mathcal{M}_j$
- ▶ **Natural** parameter: $\lambda_j \in \Omega_j$
- ▶ **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 :

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

When $\mathbf{F}_j(\lambda_j^t, \lambda_{-j}^t)$ is invertible, the update is equivalent to

$$\text{block NGVI: } \lambda_j^{t+1} = \lambda_j^t + \beta \nabla_{m_j} \mathcal{L}^t$$

Now, we give a sufficient condition when $\mathbf{F}_j(\lambda_j^t, \lambda_{-j}^t)$ is invertible.

Definition 2: Minimal Multi-linear Exp-family (MMEF)

A conditional exp-family is said to have a minimal representation when $\mathbf{m}_j(\cdot, \lambda_{-j}) : \Omega_j \rightarrow \mathcal{M}_j$ is one-to-one, $\forall \lambda_{-j} \in \Omega_{-j}$.

Theorem 2

For an MMEF representation, the FIM $\mathbf{F}_j(\lambda)$ is positive-definite and invertible for all $\lambda \in \Omega$.

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

Consider a model with a Gaussian prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})$.

$$p(\mathcal{D}, \mathbf{z}) = \prod_{n=1}^N p(\mathcal{D}_n|\mathbf{z})\mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})$$

We use the following mixture of exponential family distributions (skew-Gaussian distribution).

$$q(\mathbf{z}) = \int \mathcal{N}(\mathbf{z}|\mu + |w|\alpha, \Sigma)\mathcal{N}(w|0, 1)dw$$

The natural parameter and expectation parameter are shown below, where $c = \sqrt{2/\pi}$.

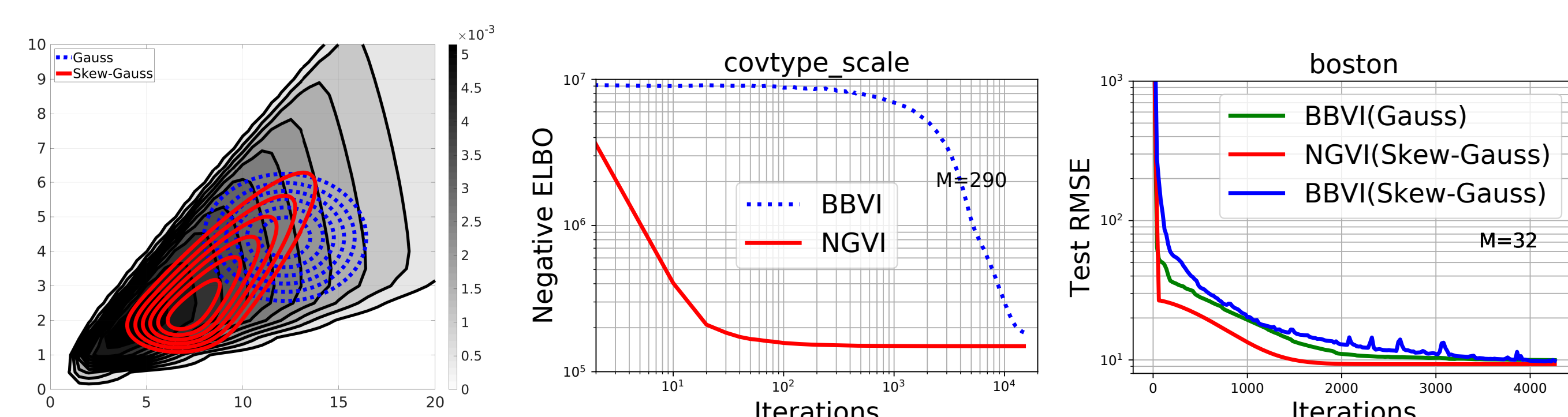
$$\lambda_{z_1} = \Sigma^{-1}\mu, \quad \lambda_{z_2} = \Sigma^{-1}\alpha, \quad \lambda_{z_3} = -\frac{1}{2}\Sigma^{-1} \\ \mathbf{m}_{z_1} = \mu + c\alpha, \quad \mathbf{m}_{z_2} = c\mu + \alpha, \quad \mathbf{m}_{z_3} = \mu\mu^T + \alpha\alpha^T + c(\mu\alpha^T + \alpha\mu^T) + \Sigma$$

The ELBO: $\mathcal{L}(\lambda) = \mathbb{E}_{q(\mathbf{z})}[\sum_{n=1}^N \log p(\mathcal{D}_n|\mathbf{z}) + \overbrace{\log \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})}^{\text{prior}} - \overbrace{\log q(\mathbf{z})}^{\text{entropy}}]$

We can re-express the update in terms of μ , Σ^{-1} , and α , where \mathbf{g}_μ^n , \mathbf{g}_α^n , and \mathbf{g}_Σ^n are defined in the paper. The derivatives about the prior and the entropy can be computed almost exactly.

$$\text{NGVI: } \Sigma^{-1} \leftarrow (1 - \beta)\Sigma^{-1} + \beta(\delta\mathbf{I} + N\mathbf{g}_\Sigma^n) \\ \mu \leftarrow \mu - \beta\Sigma(\frac{N}{1 - c^2}(\mathbf{g}_\mu^n - c\mathbf{g}_\alpha^n) + \delta\mu) \\ \alpha \leftarrow \alpha - \beta\Sigma(\frac{N}{1 - c^2}(\mathbf{g}_\alpha^n - c\mathbf{g}_\mu^n) + \delta\alpha)$$

VI using skew Gaussian distribution



Example of Mixture of Exponential Family:

Consider a model with a Gaussian prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})$.

$$p(\mathcal{D}, \mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I}) \prod_n p(\mathcal{D}_n|\mathbf{z})$$

We use a K -mixture of Gaussians shown below, where $\pi_K = 1 - \sum_{c=1}^{K-1} \pi_c$.

$$q(\mathbf{z}) = \sum_{w=1}^K \underbrace{\text{Cate}_K(w|\pi)}_{\pi_w} \mathcal{N}(\mathbf{z}|\mu_w, \Sigma_w), \text{ where } \text{Cate}_K(w|\pi) = \exp\left(\sum_{c=1}^{K-1} \mathbb{I}_c(w) \log \frac{\pi_c}{\pi_K} + \log \pi_K\right)$$

The natural parameter and expectation parameter are

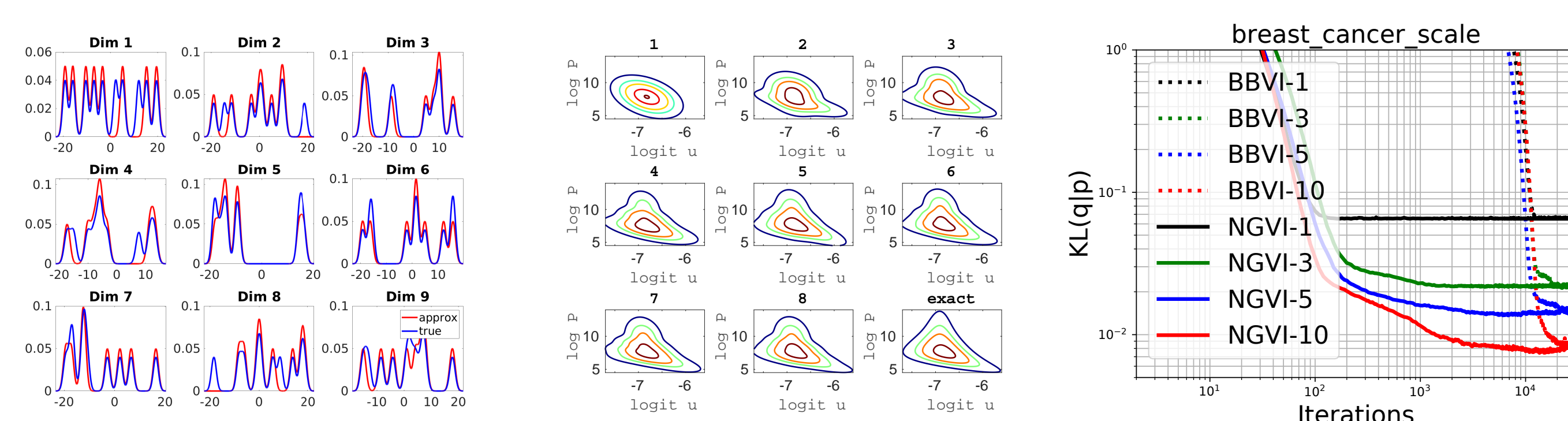
$$\lambda_z = \left\{ \Sigma_c^{-1}\mu_c, -\frac{1}{2}\Sigma_c^{-1} \right\}_{c=1}^K, \quad \mathbf{m}_z = \left\{ \pi_c\mu_c, \pi_c(\mu_c\mu_c^T + \Sigma_c) \right\}_{c=1}^K \\ \lambda_w = \left\{ \log \frac{\pi_c}{\pi_K} \right\}_{c=1}^{K-1}, \quad \mathbf{m}_w = \left\{ \pi_c \right\}_{c=1}^{K-1}$$

The ELBO: $\mathcal{L} = \mathbb{E}_{q(\mathbf{z})}[-h(\mathbf{z})]$, where $h(\mathbf{z}) := \log[q(\mathbf{z})/p(\mathbf{z})] - \sum_n \log p(\mathcal{D}_n|\mathbf{z})$.

We can re-express the update in terms of $\{\mu_c\}_{c=1}^K$, $\{\Sigma_c\}_{c=1}^K$, and $\{\pi_c\}_{c=1}^K$, where $\delta_c := \mathcal{N}(\mathbf{z}|\mu_c, \Sigma_c) / \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{z}|\mu_k, \Sigma_k)$ and \mathbf{z} is generated from $q(\mathbf{z})$.

$$\text{NGVI: } \Sigma_c^{-1} \leftarrow \Sigma_c^{-1} + \beta\delta_c [\nabla_z^2 h(\mathbf{z})] \text{ for } c = 1, \dots, K \\ \mu_c \leftarrow \mu_c - \beta\delta_c \Sigma_c [\nabla_z h(\mathbf{z})] \text{ for } c = 1, \dots, K \\ \log(\pi_c/\pi_K) \leftarrow \log(\pi_c/\pi_K) - \beta(\delta_c - \delta_K)h(\mathbf{z}) \text{ for } c = 1, \dots, K-1$$

VI using finite mixture of Gaussians.



Example of Multi-linear Exponential Family Approximation:

Consider a Bayesian model $p(\mathcal{D}, \mathbf{Z})$. We use a matrix Gaussian distribution $\mathbf{Z} \in \mathcal{R}^{d \times p}$.

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

The natural parameter and expectation parameter are

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

Using the Gauss-Newton approximation to the Hessian matrix, we obtain the following update, where the gradient is $\mathbf{G} := \nabla_{\mathbf{Z}}[-\log p(\mathcal{D}, \mathbf{Z}) + \log q(\mathbf{Z})]$ and \mathbf{Z} is sampled from $q(\mathbf{Z})$.

$$\text{block NGVI: } \mathbf{W} \leftarrow \mathbf{W} - \beta_1 \mathbf{U} \mathbf{G} \mathbf{V}, \\ \mathbf{U}^{-1} \leftarrow \mathbf{U}^{-1} + \beta_2 \mathbf{G} \mathbf{V} \mathbf{G}^T, \\ \mathbf{V}^{-1} \leftarrow \mathbf{V}^{-1} + \beta_2 \mathbf{G}^T \mathbf{U} \mathbf{G},$$

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

More Examples: For Birnbaum-Saunders distribution, exponentially modified Gaussian, Student's t, symmetric normal inverse-Gaussian, please see the appendix of the paper.

Conclusion:

We propose a new type of simple and faster natural-gradient updates for several kinds of approximations outside the existing class of exponential-family distributions.