

# Fast and Simple Natural-Gradient Variational Inference with Mixture of Exponential-family Approximations

Wu Lin†, Mohammad Emtiyaz Khan\*, Mark Schmidt†

†University of British Columbia, \*RIKEN Center for AI Project

wlin2018@cs.ubc.ca, emtiyaz.khan@riken.jp, schmidtm@cs.ubc.ca

## Abstract

Natural-gradient methods enable fast and simple algorithms for variational inference, but currently, their use is restricted to exponential-family approximations due to computational difficulties. In this paper, we extend the application of natural-gradients beyond exponential-family approximations. By proposing new types of expectation parameterizations, we derive simple updates for mixture of exponential-family approximations. Empirical results demonstrate a faster convergence of our method compared to black-box variational inference. Our work expands the scope of natural gradients for approximate Bayesian inference and makes them more widely applicable than before.

## 1. Introduction

Variational Inference provides a cheap and quick approximation to the posterior distribution and is now widely used in many areas of machine learning (Kingma and Welling, 2013; Furnston and Barber, 2010; Wainwright and Jordan, 2008; Hensman et al., 2013; Nguyen et al., 2017). In recent years, many natural-gradient methods have been proposed for variational inference (VI). However most of them are restricted to exponential-family approximations (Honkela et al., 2011; Hensman et al., 2012; Hoffman et al., 2013; Khan and Lin, 2017; Salimbeni et al., 2018; Khan et al., 2018; Zhang et al., 2018). For such approximations, the natural gradient admits a simple form and can be computed without an explicit computation of the Fisher information matrix. Unfortunately, this simplicity does not extend to other types of approximations, and derivation of such simple updates for a generic approximation remains an open problem.

In this paper, we present a new approach to obtain simple natural-gradient updates for several types of approximations outside the class of exponential-family distributions. Our main focus is the structured approximations obtained by using the mixture of exponential-family distributions. For such distributions, we propose a new type of *expectation-parameterization* that enables the derivation of simple natural-gradient updates. We apply our method to a variety of models. Our experiments establish faster convergence than black-box VI methods.

## 2. Natural-Gradient Variational Inference (NGVI)

Given a probabilistic model  $p(\mathcal{D}, \mathbf{z})$  to model the data  $\mathcal{D}$  using a latent vector  $\mathbf{z}$ , a common choice to approximate the posterior is to employ an exponential-family approximation:  $p(\mathbf{z}|\mathcal{D}) \approx q(\mathbf{z}|\boldsymbol{\lambda}_z) := h_z(\mathbf{z}) \exp[\langle \boldsymbol{\phi}_z(\mathbf{z}), \boldsymbol{\lambda}_z \rangle - A_z(\boldsymbol{\lambda}_z)]$ , where  $q$  denotes the approximating distribution,  $\boldsymbol{\lambda}_z$  is the natural parameter,  $\boldsymbol{\phi}_w(\mathbf{z})$  are the sufficient statistics,  $A_z(\boldsymbol{\lambda}_z)$  is the log-partition function, and  $h_z$  is the base measure. The approximation can be estimated

by maximizing the variational lower bound  $\mathcal{L}(\boldsymbol{\lambda}_z)$ . This optimization can be carried out by either using gradient-descent (denoted by BBVI) or using natural-gradient method (denoted by NGVI) as shown below:

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

where  $\nabla$  denotes the gradient,  $\mathbf{F}_z(\boldsymbol{\lambda}_z)$  is the Fisher information matrix (FIM) with respect to the natural parameters, and  $\alpha, \beta$  are scalar learning rates. The FIM is assumed to be invertible. In general, BBVI is computationally cheaper to perform and scales well to large data by using stochastic approximations of the gradients (Ranganath et al., 2014). However, NGVI could be simpler to compute than BBVI for some cases, e.g., in the stochastic variational inference (SVI) algorithm (Hoffman et al., 2013).

In general, a simple update can be obtained by using the expectation parameter of a *minimal* exponential family (Wainwright and Jordan, 2008) defined as  $\mathbf{m}_z(\boldsymbol{\lambda}_z) := \mathbb{E}_q[\boldsymbol{\phi}_z(\mathbf{z})] = \nabla_{\boldsymbol{\lambda}_z} A(\boldsymbol{\lambda}_z)$ . The FIM is invertible in the minimal exponential family. As shown in Khan and Nielsen (2018), the NGVI update can then be expressed as follows:

$$\text{NGVI for Exp-Family: } \boldsymbol{\lambda}_z \leftarrow \boldsymbol{\lambda}_z + \beta \nabla_{\mathbf{m}_z} \mathcal{L}(\boldsymbol{\lambda}_z), \quad (2)$$

where the gradient is computed with respect to  $\mathbf{m}_z$ . When the gradient with respect to  $\mathbf{m}_z$  is easier to compute than the gradient with respect to  $\boldsymbol{\lambda}_z$ , we get a simpler algorithm that can also potentially converge faster. This approach has been used in some recent works to derive simple natural-gradient updates (Hensman et al., 2012; Khan and Lin, 2017; Khan et al., 2018; Salimbeni et al., 2018).

Unfortunately, when the approximations are not exponential-family distributions, the NGVI update may not be written in the form (2). This is because there may not even be a notion of the expectation parameter since the distribution may not be describable by sufficient statistics. This work addresses this issue for some types of non-exponential-family approximations and derives simple natural-gradient update similar to (2).

### 3. NGVI for Mixtures of Exponential-Family Distributions

We consider the mixture distributions  $q(\mathbf{z}) := \int q(\mathbf{z}|\mathbf{w})q(\mathbf{w})d\mathbf{w}$  where  $q(\mathbf{w})$  is a mixing distribution that takes an exponential-family form as shown in (3) and  $q(\mathbf{z}|\mathbf{w})$  is the mixture component that takes the *conditional* exponential-family distribution as shown in (4).

$$q(\mathbf{w}|\boldsymbol{\lambda}_w) := h_w(\mathbf{w}) \exp[\langle \boldsymbol{\phi}_w(\mathbf{w}), \boldsymbol{\lambda}_w \rangle - A_w(\boldsymbol{\lambda}_w)], \quad (3)$$

$$q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) := h_z(\mathbf{z}, \mathbf{w}) \exp[\langle \boldsymbol{\phi}_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}_z \rangle - A_z(\boldsymbol{\lambda}_z, \mathbf{w})]. \quad (4)$$

Note that, in  $q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)$ , the sufficient statistics  $\boldsymbol{\phi}_z(\mathbf{z}, \mathbf{w})$  and log-partition  $A_z(\boldsymbol{\lambda}_z, \mathbf{w})$  depend on  $\mathbf{w}$ . However, conditioned on  $w$ , the distribution is an exponential family parametrized by natural parameter  $\boldsymbol{\lambda}_z$ . Therefore, this is a type of conditional exponential-family distribution (Xing et al., 2002; Liang et al., 2009; Lindsey, 1996; Feigin, 1981). We will give examples of this family soon, but first we define the FIM of  $q(\mathbf{z}, \mathbf{w})$ . Denoting the set of natural parameters by  $\boldsymbol{\lambda} := \{\boldsymbol{\lambda}_z, \boldsymbol{\lambda}_w\}$ , we can define the following FIM:  $\mathbf{F}_{wz}(\boldsymbol{\lambda}) := -\mathbb{E}_{q(\mathbf{z}, \mathbf{w})} [\nabla_{\boldsymbol{\lambda}}^2 \log q(\mathbf{z}, \mathbf{w})]$ . The FIM is defined by using the joint distribution of  $\mathbf{z}$  and  $\mathbf{w}$ , and it is different from the one defined over the marginal distribution  $q(\mathbf{z})$ . We also assume the minimality holds in  $q(\mathbf{w}|\boldsymbol{\lambda}_w)$  and  $q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)$  so that the FIM is invertible. The minimality is discussed in Appendix A.

We consider the natural gradient in the Riemannian metric defined by  $\mathbf{F}_{wz}(\boldsymbol{\lambda})$ . To simplify the natural-gradient update, we define the following expectation-parameters:

$$\mathbf{m}_z(\boldsymbol{\lambda}_z, \boldsymbol{\lambda}_w) := \mathbb{E}_{q(z|w)q(w)} [\boldsymbol{\phi}_z(\mathbf{z}, \mathbf{w})], \quad \mathbf{m}_w(\boldsymbol{\lambda}_w) := \mathbb{E}_{q(w)} [\boldsymbol{\phi}_w(\mathbf{w})]. \quad (5)$$

Denote the set of expectation parameters by  $\mathbf{m} := \{\mathbf{m}_z, \mathbf{m}_w\}$ . The following theorem states our main result regarding the computation of natural gradients.

**Theorem 1** *The natural-gradient update  $\boldsymbol{\lambda} \leftarrow \boldsymbol{\lambda} + \beta [F_{wz}(\boldsymbol{\lambda})]^{-1} \nabla_{\boldsymbol{\lambda}} \mathcal{L}(\boldsymbol{\lambda})$  can be performed by using the following two updates which use the gradient with respect to  $\mathbf{m}_z$  and  $\mathbf{m}_w$ :*

$$\boldsymbol{\lambda}_z \leftarrow \boldsymbol{\lambda}_z + \beta \nabla_{\mathbf{m}_z} \mathcal{L}(\boldsymbol{\lambda}), \quad \boldsymbol{\lambda}_w \leftarrow \boldsymbol{\lambda}_w + \beta \nabla_{\mathbf{m}_w} \mathcal{L}(\boldsymbol{\lambda}). \quad (6)$$

A proof-sketch is given in Appendix A. The following examples show the simplicity of the proposed update (6).

**Example 1 (Finite Mixture of Gaussians)** *The first example is regarding a finite mixture of Gaussians:  $q(\mathbf{z}) = \sum_{c=1}^K q(\mathbf{z}|w=c)q(w=c)$  where  $q(\mathbf{z}|w=c) := \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$  and  $q(w=c) = \pi_c$  with  $\sum_c \pi_c = 1$ . The natural-parameter  $\boldsymbol{\lambda}_z := \{\boldsymbol{\Sigma}_c^{-1}\boldsymbol{\mu}_c, -\frac{1}{2}\boldsymbol{\Sigma}_c^{-1}\}_{c=1}^K$  and the sufficient statistics is  $\boldsymbol{\phi}_z(\mathbf{z}, w) := \{\mathbb{I}_c(w)\mathbf{z}, \mathbb{I}_c(w)\mathbf{z}\mathbf{z}^T\}_{c=1}^K$  where  $\mathbb{I}_c(w)$  is an indicator taking a value 1 when  $w=c$  and 0 otherwise. The mixing distribution  $q(w)$  is a categorical distribution with  $\boldsymbol{\lambda}_w = \{\log \frac{\pi_c}{\pi_K}\}_{c=1}^{K-1}$  and  $\boldsymbol{\phi}_w(w) := \{\mathbb{I}_c(w)\}_{c=1}^{K-1}$ . The parameters  $\mathbf{m}_z$  for the  $c$ -th component are simply equal to the mean parameter of the Gaussian  $\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$  multiplied by  $\pi_c$ , and  $\mathbf{m}_w := \{\pi_c\}_{c=1}^{K-1}$ . The gradient can be computed using the reparameterization trick. Figure 1 demonstrates the fast convergence of our algorithm over BBVI. This result can be extended to a finite mixture of exponential-family distributions discussed in Appendix B.*

**Example 2 (Multivariate t-distribution)** *The second example is the multivariate Student's t-distribution which can be expressed as a scale mixture of Gaussian distribution:  $q(\mathbf{z}) = \int q(\mathbf{z}|w)q(w)dw = \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma}) \text{InvGam}(w|a, a)dw$ . The natural parameter and the sufficient statistics of  $q(\mathbf{z}|w)$  are  $\boldsymbol{\lambda}_z = \{\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}, -\frac{1}{2}\boldsymbol{\Sigma}^{-1}\}$  and  $\boldsymbol{\phi}_z(\mathbf{z}, w) = \{w^{-1}\mathbf{z}, w^{-1}\mathbf{z}\mathbf{z}^T\}$ . These quantities for  $q(w)$  are  $\boldsymbol{\lambda}_w := a$  and  $\boldsymbol{\phi}_w(w) := -1/w - \log w$ . The expectation parameters defined in (5) are  $\mathbf{m}_z := \{\boldsymbol{\mu}, \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma}\}$  and  $m_w := -1 - \log a + \psi(a)$ , where  $\psi$  is the digamma function. Since  $\mathbf{m}_z$  are essentially the expectation parameters of a Gaussian, the gradient with respect to them can be computed using methods described in Khan et al. (2018). The resulting update can be expressed as a perturbed Newton update. The update with respect to  $m_w$  can be easily computed using the chain rule. The full update and experimental results are included in Appendix C and G, respectively. This result can be generalized to the general class of scale mixture of Gaussian distributions. Note that the joint distribution  $\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})\text{InvGam}(w|a, a)$  indeed is a curved exponential family as shown in Appendix C. Existing methods such as Khan and Lin (2017); Khan et al. (2018) cannot be directly used.*

#### 4. Extension to Multi-Linear Exponential Family

The mixture distribution discussed above is a hierarchical distribution that contains one block of parameters at each level. We can extend our approach to a distribution with multiple blocks of parameters at a single level. Suppose we divide the vector  $\boldsymbol{\lambda} := \{\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2, \dots, \boldsymbol{\lambda}_N\}$  where  $\boldsymbol{\lambda}_j$  is the  $j$ -th block of parameters. We consider a family of distribution, which takes the following form:  $q(\mathbf{z}|\boldsymbol{\lambda}) = h_z(\mathbf{z}) \exp[f(\mathbf{z}, \boldsymbol{\lambda}) - A_z(\boldsymbol{\lambda})]$ . with  $f(\cdot)$  being a multi-linear function

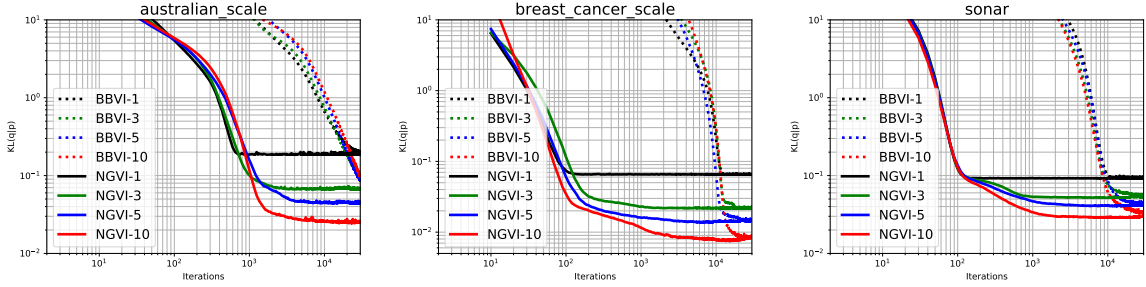


Figure 1: This figure demonstrates a fast convergence of NGVI over BBVI to approximate the posterior distribution of Bayesian logistic regression. We use a mixture of Gaussians with full covariance matrix as the approximating distribution. The number indicates the number of mixture components used. The plot shows the KL obtained using  $10^6$  MC samples, where  $p$  is the true posterior. For both algorithms, we used full-batches by using 20 MC samples to compute stochastic approximations. For BBVI, we use the Adam optimizer. We can see that our method is faster than BBVI using Adam. The main reason behind this is the simplicity of the natural-gradient update where we do not need to maintain a separate scaling vector for mean and covariance of the Gaussian, unlike Adam. This is similar to the Variational Adam method of Khan et al. (2018) and Zhang et al. (2018) which has a simpler update than Adam.

of  $\lambda_j$ , i.e., for all  $j$ , there exist functions  $\phi_j$  and  $r_j$  such that  $f$  is linear with respect to  $\lambda_j$ , i.e.,  $f(\mathbf{z}, \boldsymbol{\lambda}) := \langle \lambda_j, \phi_j(\mathbf{z}, \boldsymbol{\lambda}_{-j}) \rangle + r_j(\mathbf{z}, \boldsymbol{\lambda}_{-j})$ , where  $\boldsymbol{\lambda}_{-j}$  is the parameter vector containing all  $\boldsymbol{\lambda}$  except  $\lambda_j$ . We call this family the *multi-linear* exponential family. Obviously, an exponential family distribution parametrized by its natural parameter is a member of the multi-linear exponential family.

First, we define the natural parameter for the  $j$ -th block as  $\lambda_j$ . The expectation parameter can be defined in a similar fashion as the mixture case. The expectation parameter for the  $j$ -th block is defined as  $\mathbf{m}_j(\boldsymbol{\lambda}) := \mathbb{E}_{q(\mathbf{z})}[\phi_j(\mathbf{z}, \boldsymbol{\lambda}_{-j})]$ . By defining an FIM for the  $j$ -th block as  $\mathbf{F}_j(\boldsymbol{\lambda}) := -\mathbb{E}_{q(\mathbf{z})}[\nabla_{\lambda_j}^2 \log q(\mathbf{z}|\boldsymbol{\lambda})]$ , we can perform a *block* natural-gradient update:  $\lambda_j \leftarrow \lambda_j + \beta_j [\mathbf{F}_j(\boldsymbol{\lambda})]^{-1} \nabla_{\lambda_j} \mathcal{L}(\boldsymbol{\lambda})$ . Similarly, we assume the minimality holds for all blocks so that  $\mathbf{F}_j(\boldsymbol{\lambda})$  is invertible. The following theorem establishes an update by using the expectation parameter  $\mathbf{m}_j$ .

**Theorem 2** *The above block natural-gradient update can be performed by using the following sequential update  $\lambda_j \leftarrow \lambda_j + \beta \nabla_{\mathbf{m}_j} \mathcal{L}(\boldsymbol{\lambda})$ .*

The proof of this theorem is similar to Theorem 1. More theoretical results can be found in Appendix E and F. The following example shows the simplicity of the proposed update (6).

**Example 3 (Matrix-Variate Gaussian (MVG))** *As shown in Appendix D, the MVG distribution  $\mathcal{MN}(\mathbf{Z}|\mathbf{W}, \mathbf{U}, \mathbf{V})$  can be written in the multi-linear form. The natural-gradient update is also derived in Appendix D and below we summarize the update. We first expressing the lower bound as  $\mathcal{L}(\boldsymbol{\lambda}) = E_q[-h(\mathbf{Z})]$ . Under the Gauss-Newton approximation (Graves, 2011), the block natural-gradient update is,*

$$\mathbf{W} \leftarrow \mathbf{W} - \beta_1 \mathbf{U} \mathbf{G} \mathbf{V}, \quad \mathbf{U}^{-1} \leftarrow \mathbf{U}^{-1} + \beta_2 \mathbf{G} \mathbf{V} \mathbf{G}^\top, \quad \mathbf{V}^{-1} \leftarrow \mathbf{V}^{-1} + \beta_2 \mathbf{G}^\top \mathbf{U} \mathbf{G}, \quad (7)$$

where we sample  $\mathbf{Z}$  from the MVG distribution and evaluate the gradient  $\mathbf{G} := \nabla_{\mathbf{z}} h(\mathbf{Z})$ . These updates extend the Newton-like update obtained in Khan et al. (2018) to MVG approximations. The gradient  $\mathbf{G}$  is pre-conditioned, which is very similar to other preconditioned algorithms, such as K-FAC (Martens and Grosse, 2015; Zhang et al., 2018) and Shampoo (Gupta et al., 2018). The update can be extended to Tensor-Variate Gaussian (Ohlson et al., 2013).

## References

- Paul D Feigin. Conditional exponential families and a representation theorem for asymptotic inference. *The Annals of Statistics*, pages 597–603, 1981.
- Thomas Furnston and David Barber. Variational methods for reinforcement learning. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 241–248, 2010.
- Alex Graves. Practical variational inference for neural networks. In *Advances in neural information processing systems*, pages 2348–2356, 2011.
- Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor optimization. In *Proceedings of the 35th International Conference on Machine Learning*, pages 1842–1850, 2018.
- James Hensman, Magnus Rattray, and Neil D Lawrence. Fast variational inference in the conjugate exponential family. In *Advances in neural information processing systems*, pages 2888–2896, 2012.
- James Hensman, Nicolo Fusi, and Neil D Lawrence. Gaussian processes for big data. *arXiv preprint arXiv:1309.6835*, 2013.
- Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303–1347, 2013.
- A. Honkela, T. Raiko, M. Kuusela, M. Torniio, and J. Karhunen. Approximate Riemannian conjugate gradient learning for fixed-form variational Bayes. 11:3235–3268, 2011.
- Mohammad Khan and Wu Lin. Conjugate-computation variational inference: Converting variational inference in non-conjugate models to inferences in conjugate models. In *Artificial Intelligence and Statistics*, pages 878–887, 2017.
- Mohammad Emtiyaz Khan and Didrik Nielsen. Fast yet simple natural-gradient descent for variational inference in complex models. *arXiv preprint arXiv:1807.04489*, 2018.
- Mohammad Emtiyaz Khan, Didrik Nielsen, Voot Tangkaratt, Wu Lin, Yarin Gal, and Akash Srivastava. Fast and scalable Bayesian deep learning by weight-perturbation in Adam. In *Proceedings of the 35th International Conference on Machine Learning*, pages 2611–2620, 2018.
- Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Erich L Lehmann and George Casella. *Theory of point estimation*. Springer Science & Business Media, 2006.
- Percy Liang, Michael I Jordan, and Dan Klein. Learning from measurements in exponential families. In *Proceedings of the 26th annual international conference on machine learning*, pages 641–648. ACM, 2009.

- James K Lindsey. *Parametric statistical inference*. Oxford University Press, 1996.
- James Martens and Roger Grosse. Optimizing neural networks with Kronecker-factored approximate curvature. In *International Conference on Machine Learning*, pages 2408–2417, 2015.
- Cuong V Nguyen, Yingzhen Li, Thang D Bui, and Richard E Turner. Variational continual learning. *arXiv preprint arXiv:1710.10628*, 2017.
- Martin Ohlson, M Rauf Ahmad, and Dietrich Von Rosen. The multilinear normal distribution: Introduction and some basic properties. *Journal of Multivariate Analysis*, 113:37–47, 2013.
- Manfred Opper and Cédric Archambeau. The variational gaussian approximation revisited. *Neural computation*, 21(3):786–792, 2009.
- Rajesh Ranganath, Sean Gerrish, and David Blei. Black box variational inference. In *Artificial Intelligence and Statistics*, pages 814–822, 2014.
- Hugh Salimbeni, Stefanos Eleftheriadis, and James Hensman. Natural gradients in practice: Non-conjugate variational inference in gaussian process models. 2018.
- M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1–2:1–305, 2008.
- Eric P Xing, Michael I Jordan, and Stuart Russell. A generalized mean field algorithm for variational inference in exponential families. In *Proceedings of the Nineteenth conference on Uncertainty in Artificial Intelligence*, pages 583–591. Morgan Kaufmann Publishers Inc., 2002.
- G. Zhang, S. Sun, D. Duvenaud, and R. Grosse. Noisy natural gradient as variational inference. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.

## Appendix A. Proof of Theorem 1

The NGVI update of (6) updates the parameters  $\boldsymbol{\lambda}_z$  and  $\boldsymbol{\lambda}_w$  separately. This is possible because the FIM is in fact block-diagonal. The following lemma states this result.

**Lemma 3** *The FIM  $\mathbf{F}_{wz}(\boldsymbol{\lambda})$  is block-diagonal with two blocks:*

$$\mathbf{F}_{wz}(\boldsymbol{\lambda}) = \begin{bmatrix} \mathbf{F}_z(\boldsymbol{\lambda}) & \mathbf{0} \\ \mathbf{0} & \mathbf{F}_w(\boldsymbol{\lambda}_w) \end{bmatrix}, \quad (8)$$

with the blocks defined as follows:

$$\mathbf{F}_z(\boldsymbol{\lambda}) := -\mathbb{E}_{q(z,w)} [\nabla_{\lambda_z}^2 \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)], \quad \mathbf{F}_w(\boldsymbol{\lambda}_w) := -\mathbb{E}_{q(w)} [\nabla_{\lambda_w}^2 \log q(\mathbf{w}|\boldsymbol{\lambda}_w)] \quad (9)$$

**Proof** By definition of FIM, we have

$$\mathbf{F}_{wz}(\boldsymbol{\lambda}) = -\mathbb{E}_{q(z,w|\boldsymbol{\lambda})} [\nabla_{\boldsymbol{\lambda}}^2 \log q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda})] \quad (10)$$

$$= -\mathbb{E}_{q(z,w|\boldsymbol{\lambda})} \begin{bmatrix} \nabla_{\lambda_z}^2 \log q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda}) & \nabla_{\lambda_w} \nabla_{\lambda_z} \log q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda}) \\ \nabla_{\lambda_z} \nabla_{\lambda_w} \log q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda}) & \nabla_{\lambda_w}^2 \log q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda}) \end{bmatrix} \quad (11)$$

$$= -\mathbb{E}_{q(z,w|\boldsymbol{\lambda})} \begin{bmatrix} \nabla_{\lambda_z}^2 (\log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) + \log q(\mathbf{w}|\boldsymbol{\lambda}_w)) & \nabla_{\lambda_w} \nabla_{\lambda_z} (\log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) + \log q(\mathbf{w}|\boldsymbol{\lambda}_w)) \\ \nabla_{\lambda_z} \nabla_{\lambda_w} (\log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) + \log q(\mathbf{w}|\boldsymbol{\lambda}_w)) & \nabla_{\lambda_w}^2 (\log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) + \log q(\mathbf{w}|\boldsymbol{\lambda}_w)) \end{bmatrix} \quad (12)$$

$$= -\mathbb{E}_{q(z,w|\boldsymbol{\lambda})} \begin{bmatrix} \nabla_{\lambda_z}^2 \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) & \mathbf{0} \\ \mathbf{0} & \nabla_{\lambda_w}^2 \log q(\mathbf{w}|\boldsymbol{\lambda}_w) \end{bmatrix} \quad (13)$$

$$= - \begin{bmatrix} \mathbb{E}_{q(z,w|\boldsymbol{\lambda})} [\nabla_{\lambda_z}^2 \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)] & \mathbf{0} \\ \mathbf{0} & \mathbb{E}_{q(w|\boldsymbol{\lambda}_w)} [\nabla_{\lambda_w}^2 \log q(\mathbf{w}|\boldsymbol{\lambda}_w)] \end{bmatrix} \quad (14)$$

$$= \begin{bmatrix} \mathbf{F}_z(\boldsymbol{\lambda}) & \mathbf{0} \\ \mathbf{0} & \mathbf{F}_w(\boldsymbol{\lambda}_w) \end{bmatrix}, \quad (15)$$

where we can move from step 12 to step 13 since  $\boldsymbol{\lambda}_z$  and  $\boldsymbol{\lambda}_w$  are not tied.  $\blacksquare$

By Lemma 3, the updates can therefore be carried out separately. The updates for  $\boldsymbol{\lambda}_w$  can be expressed in terms of its mean parameters by using the results of Khan and Nielsen (2018) because  $q(\mathbf{w})$  is an exponential family distribution. Therefore, we only need to derive the update for  $\boldsymbol{\lambda}_z$ . The following theorem states that the corresponding FIM can be written as the derivative of  $\mathbf{m}_z$  with respect to  $\boldsymbol{\lambda}_z$ .

**Lemma 4** *The FIM matrix  $\mathbf{F}_z$  is equal to the derivative of the expectation parameter  $\mathbf{m}_z(\boldsymbol{\lambda})$ :*

$$\mathbf{F}_z(\boldsymbol{\lambda}) := \nabla_{\lambda_z} \mathbf{m}_z(\boldsymbol{\lambda}) \quad (16)$$

**Proof** Recall that

$$\mathbb{E}_{q(z,w|\boldsymbol{\lambda})} [\nabla_{\lambda_z} \log q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda})] = \mathbb{E}_{q(z,w|\boldsymbol{\lambda})} [\nabla_{\lambda_z} \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)] \quad (17)$$

$$= \mathbb{E}_{q(z,w|\boldsymbol{\lambda})} [[\nabla_{\lambda_z} q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)] / q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)] \quad (18)$$

$$= \mathbb{E}_{q(w|\boldsymbol{\lambda}_w)} [\nabla_{\lambda_z} q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)] \quad (19)$$

$$= \nabla_{\lambda_z} \mathbb{E}_{q(z,w|\boldsymbol{\lambda})} [1] \quad (20)$$

$$= \mathbf{0} \quad (21)$$

Since  $q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) = h_z(\mathbf{z}, \mathbf{w}) \exp[\langle \boldsymbol{\phi}_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}_z \rangle - A_z(\boldsymbol{\lambda}_z, \mathbf{w})]$ , we have

$$\mathbf{0} = \mathbb{E}_{q(z, \mathbf{w}|\boldsymbol{\lambda})} [\nabla_{\lambda_z} \log q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda})] \quad (22)$$

$$= \mathbb{E}_{q(z, \mathbf{w}|\boldsymbol{\lambda})} [\nabla_{\lambda_z} \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)] \quad (23)$$

$$= \mathbb{E}_{q(z, \mathbf{w}|\boldsymbol{\lambda})} [\nabla_{\lambda_z} (\log h_z(\mathbf{z}, \mathbf{w}) + \langle \boldsymbol{\phi}_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}_z \rangle - A_z(\boldsymbol{\lambda}_z, \mathbf{w}))] \quad (24)$$

$$= \mathbb{E}_{q(z, \mathbf{w}|\boldsymbol{\lambda})} [\nabla_{\lambda_z} (\langle \boldsymbol{\phi}_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}_z \rangle - A_z(\boldsymbol{\lambda}_z, \mathbf{w}))] \quad (25)$$

$$= \mathbb{E}_{q(z, \mathbf{w}|\boldsymbol{\lambda})} [\boldsymbol{\phi}_z(\mathbf{z}, \mathbf{w})] - \mathbb{E}_{q(z, \mathbf{w}|\boldsymbol{\lambda})} [\nabla_{\lambda_z} A_z(\boldsymbol{\lambda}_z, \mathbf{w})] \quad (26)$$

$$= \mathbf{m}_z - \mathbb{E}_{q(z, \mathbf{w}|\boldsymbol{\lambda})} [\nabla_{\lambda_z} A_z(\boldsymbol{\lambda}_z, \mathbf{w})] \quad (27)$$

$$= \mathbf{m}_z - \mathbb{E}_{q(\mathbf{w}|\boldsymbol{\lambda}_w)} [\nabla_{\lambda_z} A_z(\boldsymbol{\lambda}_z, \mathbf{w})] \quad (28)$$

Therefore, we have

$$\mathbf{F}_z(\boldsymbol{\lambda}) = -\mathbb{E}_{q(z, \mathbf{w})} [\nabla_{\lambda_z}^2 \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)] \quad (29)$$

$$= -\mathbb{E}_{q(z, \mathbf{w}|\boldsymbol{\lambda})} [\nabla_{\lambda_z}^2 (\log h_z(\mathbf{z}, \mathbf{w}) + \langle \boldsymbol{\phi}_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}_z \rangle - A_z(\boldsymbol{\lambda}_z, \mathbf{w}))] \quad (30)$$

$$= \mathbb{E}_{q(z, \mathbf{w}|\boldsymbol{\lambda})} [\nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w})] \quad (31)$$

$$= \mathbb{E}_{q(\mathbf{w}|\boldsymbol{\lambda}_w)} [\nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w})] \quad (32)$$

$$= \nabla_{\lambda_z} \mathbb{E}_{q(\mathbf{w}|\boldsymbol{\lambda}_w)} [\nabla_{\lambda_z} A_z(\boldsymbol{\lambda}_z, \mathbf{w})] \quad (33)$$

$$= \nabla_{\lambda_z} \mathbf{m}_z \quad (34)$$

■

Using this result and applying chain rule, we can write the following:

$$\mathbf{F}_z(\boldsymbol{\lambda}) [\nabla_{m_z} \mathcal{L}(\boldsymbol{\lambda})] = \nabla_{\lambda_z} \mathbf{m}_z(\boldsymbol{\lambda}) [\nabla_{m_z} \mathcal{L}(\boldsymbol{\lambda})] = \nabla_{\lambda_z} \mathcal{L}(\boldsymbol{\lambda}) \quad (35)$$

If the FIM is invertible, we can express the natural-gradient with respect to  $\boldsymbol{\lambda}_z$  as the gradient with respect to  $\mathbf{m}_z$ . The following lemma establishes the conditions under which the FIM is invertible. This condition is just a sufficient condition. There are cases where FIM is invertible but the condition does not hold.

**Lemma 5** *Let's denote the probability space of  $q(\mathbf{w})$  by  $\Omega$ . If there exists an area  $S \subseteq \Omega$  such that  $S$  is a non-zero measure and  $q(\mathbf{z}|\mathbf{w})$  is a minimal exponential-family distribution for all  $\mathbf{w} \in S$ , then  $\mathbf{F}_z(\boldsymbol{\lambda})$  is invertible.*

**Proof** Conditioned on  $\mathbf{w}$ , we know that  $q(\mathbf{z}|\mathbf{w})$  is an exponential-family distribution as shown below.

$$q(\mathbf{z}|\mathbf{w}) = h_z(\mathbf{z}, \mathbf{w}) \exp[\langle \boldsymbol{\phi}_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}_z \rangle - A_z(\boldsymbol{\lambda}_z, \mathbf{w})] \quad (36)$$

Note that

$$\mathbb{E}_{q(z|\mathbf{w})} [\nabla_{\lambda_z}^2 \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)] \quad (37)$$

$$= \mathbb{E}_{q(z|\mathbf{w})} [\nabla_{\lambda_z} [\boldsymbol{\phi}_z(\mathbf{z}, \mathbf{w}) - \nabla_{\lambda_z} A_z(\boldsymbol{\lambda}_z, \mathbf{w})]] \quad (38)$$

$$= \mathbb{E}_{q(z|\mathbf{w})} [-\nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w})] \quad (39)$$

$$= -\nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w}) \quad (40)$$



Let's denote  $\mathbf{e}_z = \frac{\nabla_{\lambda_z} q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)}{q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)}$ . Conditioned on  $\mathbf{w}$ , we have

$$\nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w}) \quad (41)$$

$$= -\mathbb{E}_{q(z|\mathbf{w})} \left[ \nabla_{\lambda_z}^2 \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) \right] \quad (42)$$

$$= -\mathbb{E}_{q(z|\mathbf{w})} \left[ \nabla_{\lambda_z} \left( \frac{\nabla_{\lambda_z} q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)}{q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)} \right) \right] \quad (43)$$

$$= -\mathbb{E}_{q(z|\mathbf{w})} \left[ \frac{\nabla_{\lambda_z}^2 q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)}{q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)} - \mathbf{e}_z \mathbf{e}_z^T \right] \quad (44)$$

$$= \mathbb{E}_{q(z|\mathbf{w})} \left[ -\frac{\nabla_{\lambda_z}^2 q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)}{q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)} \right] + \mathbb{E}_{q(z|\mathbf{w})} [\mathbf{e}_z \mathbf{e}_z^T] \quad (45)$$

$$= \underbrace{\nabla_{\lambda_z}^2 \mathbb{E}_{q(z|\mathbf{w})} [-1]}_{\mathbf{0}} + \mathbb{E}_{q(z|\mathbf{w})} [\mathbf{e}_z \mathbf{e}_z^T] \quad (46)$$

$$= \mathbb{E}_{q(z|\mathbf{w})} [\mathbf{e}_z \mathbf{e}_z^T] \succeq \mathbf{0}, \quad (47)$$

where we use the chain rule to move from step 43 to step 44.

Since  $q(\mathbf{z}|\mathbf{w})$  is a minimal exponential-family distribution for all  $\mathbf{w} \in S$ , by Proposition 3.1 of [Wainwright and Jordan \(2008\)](#), we know that  $A_z(\boldsymbol{\lambda}_z, \mathbf{w})$  is strictly convex with respect to  $\boldsymbol{\lambda}_z$  given  $\mathbf{w} \in S$  is known, which implies that  $\nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w}) \succ \mathbf{0}$ .

Using Eq. 32, we have

$$\mathbf{F}_z(\boldsymbol{\lambda}) = -\mathbb{E}_{q(z, \mathbf{w})} \left[ \nabla_{\lambda_z}^2 \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) \right] \quad (48)$$

$$= \mathbb{E}_{q(\mathbf{w}|\lambda_w)} \left[ \nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w}) \right] \quad (49)$$

Given any nonzero vector  $\mathbf{a}$ , we have

$$\mathbf{a}^T \mathbf{F}_z(\boldsymbol{\lambda}) \mathbf{a} \quad (50)$$

$$= \mathbf{a}^T \left[ \int_{\mathbf{w} \in S} q(\mathbf{w}|\lambda_w) \nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w}) d\mathbf{w} + \int_{\mathbf{w} \in (\Omega-S)} q(\mathbf{w}|\lambda_w) \nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w}) d\mathbf{w} \right] \mathbf{a} \quad (51)$$

$$= \int_{\mathbf{w} \in S} q(\mathbf{w}|\lambda_w) [\mathbf{a}^T (\nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w})) \mathbf{a}] d\mathbf{w} + \int_{\mathbf{w} \in (\Omega-S)} q(\mathbf{w}|\lambda_w) \underbrace{[\mathbf{a}^T (\nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w})) \mathbf{a}]}_{\geq 0} d\mathbf{w} \quad (52)$$

$$\geq \int_{\mathbf{w} \in S} q(\mathbf{w}|\lambda_w) \underbrace{[\mathbf{a}^T (\nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w})) \mathbf{a}]}_{> 0} d\mathbf{w} \quad (53)$$

$$> 0. \quad (54)$$

We use Eq 47 to move from step 52 to step 53. Furthermore, we use the strict convexity of  $A_z(\boldsymbol{\lambda}_z, \mathbf{w})$  with respect to  $\boldsymbol{\lambda}_z$  when  $\mathbf{w} \in S$  and the nonzero measure  $S$  to move from step 53 to step 54.

Therefore, we know that  $\mathbf{F}_z(\boldsymbol{\lambda}) \succ \mathbf{0}$  and  $\mathbf{F}_z(\boldsymbol{\lambda})$  is invertible.  $\blacksquare$

## Appendix B. Updates for Finite Mixture of Gaussian Distributions

Let's consider  $q(w) = \text{Cate}(w|\boldsymbol{\pi})$  and  $q(\mathbf{z}|w) = \text{ExpFmy}(\mathbf{z}|\boldsymbol{\lambda}_w)$ . We assume  $q(\mathbf{z}|w)$  is (implicitly) re-parameterizable.

**Lemma 6**  $q(\mathbf{z}|w)$  is a conditional exponential family distribution.

**Proof**

$$q(\mathbf{z}|w) = \sum_{c=1}^K \mathbb{I}_c(w) \text{ExpFmy}(\mathbf{z}|\boldsymbol{\lambda}_{z_c}) \quad (55)$$

$$= \sum_{c=1}^K \mathbb{I}_c(w) h_z(\mathbf{z}) \exp[\boldsymbol{\lambda}_{z_c}^T \boldsymbol{\phi}_z(\mathbf{z}) - A_z(\boldsymbol{\lambda}_{z_c})] \quad (56)$$

$$= h_z(\mathbf{z}) \exp\left\{ \sum_{c=1}^K (\mathbb{I}_c(w) \boldsymbol{\phi}_z(\mathbf{z}), \boldsymbol{\lambda}_{z_c}) - \sum_{c'=1}^K \mathbb{I}_{c'}(w) A_z(\boldsymbol{\lambda}_{z_{c'}}) \right\} \quad (57)$$

■

Firstly, we give the update for  $q(\mathbf{z}|w)$ . The natural parameter and sufficient statistics  $\{\boldsymbol{\lambda}_{z_c}\}_{c=1}^K$  and  $\{\mathbb{I}_c(w) \boldsymbol{\phi}_z(\mathbf{z})\}_{c=1}^K$ .

The update is

$$\boldsymbol{\lambda}_{z_c} \leftarrow \boldsymbol{\lambda}_{z_c} + \beta_z \nabla_{m_{z_c}} \mathcal{L}(\boldsymbol{\lambda}), \quad (58)$$

where  $m_{z_c} = \mathbb{E}_{q(w,z)} [\mathbb{I}_c(w) \boldsymbol{\phi}_z(\mathbf{z})]$ .

We consider the following model.

$$p(\mathcal{D}, \mathbf{z}) = \prod_{n=1}^N p(\mathcal{D}_n | \mathbf{z}) p(\mathbf{z}) \quad (59)$$

Let denote  $h_n(\mathbf{z}) := -(\log p(\mathcal{D}_n | \mathbf{z}) + \log p(\mathbf{z})/N)$ . The lower bound can be expressed as

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(w,z)} \left[ - \sum_{n=1}^N h_n(\mathbf{z}) - \log q(\mathbf{z}, w) \right] \quad (60)$$

$$= \mathbb{E}_{q(w,z)} \left[ - \sum_{n=1}^N h_n(\mathbf{z}) - \log \sum_{c'=1}^K \pi_{c'} \text{ExpFmy}(\mathbf{z}|\boldsymbol{\lambda}_{z_{c'}}) \right]. \quad (61)$$

In general,  $\nabla_{m_{z_c}} \mathcal{L}(\boldsymbol{\lambda})$  can be computed

$$\nabla_{m_{z_c}} \mathcal{L}(\boldsymbol{\lambda}) = (\nabla_{\boldsymbol{\lambda}_{z_c}} \mathbf{m}_{z_c})^{-1} \nabla_{\boldsymbol{\lambda}_{z_c}} \mathcal{L}(\boldsymbol{\lambda}) \quad (62)$$

$$= (\nabla_{\boldsymbol{\lambda}_{z_c}} \mathbb{E}_{q(w,z)} [\mathbb{I}_c(w) \boldsymbol{\phi}_z(\mathbf{z})])^{-1} \nabla_{\boldsymbol{\lambda}_{z_c}} \mathcal{L}(\boldsymbol{\lambda}), \quad (63)$$

where

$$\nabla_{\boldsymbol{\lambda}_{z_c}} \mathbb{E}_{q(w,z)} [\mathbb{I}_c(w) \boldsymbol{\phi}_z(\mathbf{z})] = \nabla_{\boldsymbol{\lambda}_{z_c}} \int \pi_c q(\mathbf{z}|w=c) \boldsymbol{\phi}_z(\mathbf{z}) d\mathbf{z} \quad (64)$$

$$= \int \pi_c q(\mathbf{z}|w=c) \nabla_z [\boldsymbol{\phi}_z(\mathbf{z})] [\nabla_{\boldsymbol{\lambda}_{z_c}} \mathbf{z}] d\mathbf{z} \quad (65)$$

$\nabla_{\lambda_{z_c}} \mathcal{L}(\boldsymbol{\lambda})$  can be computed as below.

$$\nabla_{\lambda_{z_c}} \mathcal{L}(\boldsymbol{\lambda}) = \nabla_{\lambda_{z_c}} \mathbb{E}_{q(\mathbf{z})} \left[ - \sum_{n=1}^N h_n(\mathbf{z}) \right] - \nabla_{\lambda_{z_c}} \mathbb{E}_{q(\mathbf{z})} \left[ \log \sum_{c'=1}^K \pi_{c'} \text{ExpFmy}(\mathbf{z} | \boldsymbol{\lambda}_{z_{c'}}) \right] \quad (66)$$

$$= \int \pi_c \nabla_{\lambda_{z_c}} q(\mathbf{z} | w = c) \left[ - \sum_{n=1}^N h_n(\mathbf{z}) - \log \sum_{c'=1}^K \pi_{c'} \text{ExpFmy}(\mathbf{z} | \boldsymbol{\lambda}_{z_{c'}}) \right] d\mathbf{z} + \underbrace{\int q(\mathbf{z}) [\nabla_{\lambda_{z_c}} \log q(\mathbf{z})] d\mathbf{z}}_0 \quad (67)$$

$$= \int \pi_c q(\mathbf{z} | w = c) \left[ \nabla_z \left( - \sum_{n=1}^N h_n(\mathbf{z}) - \log \sum_{c'=1}^K \pi_{c'} \text{ExpFmy}(\mathbf{z} | \boldsymbol{\lambda}_{z_{c'}}) \right) \right] [\nabla_{\lambda_{z_c}} \mathbf{z}] d\mathbf{z} \quad (68)$$

Now, we give the update for  $q(w | \boldsymbol{\lambda}_w) = \text{Cate}(w | \boldsymbol{\pi})$ , which is an exponential family distribution as shown below.

$$\text{Cate}(w | \boldsymbol{\pi}) = \exp \left\{ \left( \sum_{c=1}^{K-1} \mathbb{I}_c(w) \log \frac{\pi_c}{\pi_K} \right) + \log \pi_K \right\}, \quad (69)$$

where  $\pi_K = 1 - \sum_{c=1}^{K-1} \pi_c$

Its natural parameter and expectation parameter are

$$\boldsymbol{\lambda}_w = \left\{ \log \frac{\pi_c}{\pi_K} \right\}_{c=1}^{K-1} \quad (70)$$

$$\mathbf{m}_w = \left\{ \mathbb{E}_{q(w)} [\mathbb{I}_c(w)] \right\}_{c=1}^{K-1} = \left\{ \pi_c \right\}_{c=1}^{K-1} \quad (71)$$

Recall that  $\mathcal{L}(\boldsymbol{\lambda})$  can be expressed as

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(w, \mathbf{z})} \left[ - \sum_{n=1}^N h_n(\mathbf{z}) - \log \sum_{c'=1}^K \pi_{c'} \text{ExpFmy}(\mathbf{z} | \boldsymbol{\lambda}_{z_{c'}}) \right]. \quad (72)$$

$\nabla_{\pi_c} \mathcal{L}(\boldsymbol{\lambda})$  can be computed as below.

$$\nabla_{\pi_c} \mathcal{L}(\boldsymbol{\lambda}) = \nabla_{\pi_c} \mathbb{E}_{q(\mathbf{z})} \left[ - \sum_{n=1}^N h_n(\mathbf{z}) \right] - \nabla_{\pi_c} \mathbb{E}_{q(\mathbf{z})} \left[ \log \sum_{c'=1}^K \pi_{c'} \text{ExpFmy}(\mathbf{z} | \boldsymbol{\lambda}_{z_{c'}}) \right] \quad (73)$$

$$= \int q(\mathbf{z} | w = c) \left[ - \sum_{n=1}^N h_n(\mathbf{z}) - \log \sum_{c'=1}^K \pi_{c'} \text{ExpFmy}(\mathbf{z} | \boldsymbol{\lambda}_{z_{c'}}) \right] d\mathbf{z} + \underbrace{\int q(\mathbf{z}) [\nabla_{\pi_c} \log q(\mathbf{z})] d\mathbf{z}}_0 \quad (74)$$

$$= \int q(\mathbf{z} | w = c) \left[ - \sum_{n=1}^N h_n(\mathbf{z}) - \log \sum_{c'=1}^K \pi_{c'} \text{ExpFmy}(\mathbf{z} | \boldsymbol{\lambda}_{z_{c'}}) \right] d\mathbf{z} \quad (75)$$

Observe that in Eq 65, 68, 75, we have to use at least  $K$  MC samples to compute the gradients for a  $K$ -mixture distribution. To address this, we propose to use importance

sampling to reduce the number of MC samples. The gradients are computed as below.

$$\nabla_{\lambda_{z_c}} \mathbb{E}_{q(w,z)} [\mathbb{I}(w=c) \phi_z(\mathbf{z})] = \int q(\mathbf{z}) \frac{\pi_c q(\mathbf{z}|w=c)}{q(\mathbf{z})} [\nabla_z \phi_z(\mathbf{z})] [\nabla_{\lambda_{z_c}} \mathbf{z}] d\mathbf{z} \quad (76)$$

$$= \mathbb{E}_{q(z)} [q(w=c|\mathbf{z}) \nabla_z \phi_z(\mathbf{z}) \nabla_{\lambda_{z_c}} \mathbf{z}] \quad (77)$$

$$\nabla_{\lambda_{z_c}} \mathcal{L}(\boldsymbol{\lambda}) = \int q(\mathbf{z}) \frac{\pi_c q(\mathbf{z}|w=c)}{q(\mathbf{z})} \nabla_z \left[ -\sum_{n=1}^N h_n(\mathbf{z}) - \log \sum_{c'=1}^K \pi_{c'} \text{ExpFmy}(\mathbf{z}|\boldsymbol{\lambda}_{z_{c'}}) \right] [\nabla_{\lambda_{z_c}} \mathbf{z}] d\mathbf{z} \quad (78)$$

$$= \mathbb{E}_{q(z)} \left[ q(w=c|\mathbf{z}) \nabla_z \left( -\sum_{n=1}^N h_n(\mathbf{z}) - \log \sum_{c'=1}^K \pi_{c'} \text{ExpFmy}(\mathbf{z}|\boldsymbol{\lambda}_{z_{c'}}) \right) \nabla_{\lambda_{z_c}} \mathbf{z} \right], \quad (79)$$

$$\nabla_{\pi_c} \mathcal{L}(\boldsymbol{\lambda}) = \int q(\mathbf{z}) \frac{q(\mathbf{z}|w=c)}{q(\mathbf{z})} \left[ -\sum_{n=1}^N h_n(\mathbf{z}) - \log \sum_{c'=1}^K \pi_{c'} \text{ExpFmy}(\mathbf{z}|\boldsymbol{\lambda}_{z_{c'}}) \right] d\mathbf{z} \quad (80)$$

$$= \mathbb{E}_{q(z)} \left[ \frac{q(\mathbf{z}|w=c)}{q(\mathbf{z})} \left( -\sum_{n=1}^N h_n(\mathbf{z}) - \log \sum_{c'=1}^K \pi_{c'} \text{ExpFmy}(\mathbf{z}|\boldsymbol{\lambda}_{z_{c'}}) \right) \right] \quad (81)$$

In the mixture of Gaussian case, we can directly compute the gradient  $\nabla_{m_{z_c}} \mathcal{L}(\boldsymbol{\lambda})$ . Recall that the natural parameters of  $q(\mathbf{z}|w)$  are  $\lambda_z = \{\boldsymbol{\Sigma}_c^{-1} \boldsymbol{\mu}_c, -\frac{1}{2} \boldsymbol{\Sigma}_c^{-1}\}_{c=1}^K$ . The expectation parameters of  $q(\mathbf{z}|w)$  are  $m_z = \{\pi_c \boldsymbol{\mu}_c, \pi_c (\boldsymbol{\mu}_c \boldsymbol{\mu}_c^T + \boldsymbol{\Sigma}_c)\}_{c=1}^K$ .

By the chain rule, we know that

$$\nabla_{m_{z_c^1}} \mathcal{L}(\boldsymbol{\lambda}) = \frac{1}{\pi_c} (\nabla_{\mu_c} \mathcal{L}(\boldsymbol{\lambda}) - 2 \nabla_{\Sigma_c} \mathcal{L}(\boldsymbol{\lambda}) \boldsymbol{\mu}_c) \quad (82)$$

$$\nabla_{m_{z_c^2}} \mathcal{L}(\boldsymbol{\lambda}) = \frac{1}{\pi_c} (\nabla_{\Sigma_c} \mathcal{L}(\boldsymbol{\lambda})) \quad (83)$$

By the Bonnet's and Price's theorem, we have

$$\nabla_{\mu_c} \mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(z)} \left[ q(w=c|\mathbf{z}) \nabla_z \left( -\sum_{n=1}^N h_n(\mathbf{z}) - \log \sum_{c'=1}^K \pi_{c'} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{c'}, \boldsymbol{\Sigma}_{c'}) \right) \right] \quad (84)$$

$$\nabla_{\Sigma_c} \mathcal{L}(\boldsymbol{\lambda}) = \frac{1}{2} \mathbb{E}_{q(z)} \left[ q(w=c|\mathbf{z}) \nabla_z^2 \left( -\sum_{n=1}^N h_n(\mathbf{z}) - \log \sum_{c'=1}^K \pi_{c'} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{c'}, \boldsymbol{\Sigma}_{c'}) \right) \right]. \quad (85)$$

### Appendix C. Update for Multivariate t-Distribution

In this section, we derive the update for the multivariate t-distribution. The update can be extended to the class of normal variance-mean mixture including the scale mixture of Gaussians.

**Lemma 7**  $\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma}) \text{InvGam}(w|a, a)$ , where  $\mathbf{z} \in \mathcal{R}^d$ , is a curved exponential family distribution (Lehmann and Casella, 2006).

**Proof** Let  $\Lambda_1 = -\frac{1}{2}\Sigma^{-1}$ ,  $\lambda_2 = \Sigma^{-1}\mu$ ,  $\lambda_3 = -\frac{1}{2}\mu^T\Sigma^{-1}\mu$ ,  $\lambda_4 = -a$ . The distribution on  $\mathbf{z} \in \mathcal{R}^d$  and  $w \in \mathcal{R}_{++}$  can be expressed as follows.

$$\mathcal{N}(\mathbf{z}|\mu, w\Sigma)\text{InvGam}(w|a, a) \quad (86)$$

$$= \det(2\pi w\Sigma)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\mathbf{z} - \mu)^T (w\Sigma)^{-1} (\mathbf{z} - \mu)\right\} \frac{a^a}{\Gamma(a)} w^{-a-1} \exp\left\{-\frac{a}{w}\right\} \quad (87)$$

$$= (2\pi w)^{-d/2} w^{-1} \exp\left\{-\frac{1}{2}(\mathbf{z} - \mu)^T (w\Sigma)^{-1} (\mathbf{z} - \mu) - \frac{1}{2} \log \det(\Sigma) - \frac{a}{w} - a \log w - (\log \Gamma(a) - a \log a)\right\} \quad (88)$$

$$= (2\pi w)^{-d/2} w^{-1} \exp\left\{\langle -\frac{1}{2}\Sigma^{-1}, w\mathbf{z}\mathbf{z}^T \rangle + \langle \Sigma^{-1}\mu, w\mathbf{z} \rangle + \langle -\frac{1}{2}\mu^T\Sigma^{-1}\mu, w^{-1} \rangle + \langle -a, w^{-1} + \log w \rangle - [\log \Gamma(a) - a \log a + \frac{1}{2} \log \det(\Sigma)]\right\} \quad (89)$$

$$= (2\pi w)^{-d/2} w^{-1} \exp\left\{\langle \Lambda_1, w\mathbf{z}\mathbf{z}^T \rangle + \langle \lambda_2, w\mathbf{z} \rangle + \langle \lambda_3, w^{-1} \rangle + \langle \lambda_4, w^{-1} + \log w \rangle - [\log \Gamma(-\lambda_4) + \lambda_4 \log(-\lambda_4) - \frac{1}{2} \log \det(-2\Lambda_1)]\right\}, \quad (90)$$

It is easy to see that the joint distribution is an exponential family distribution. Note that  $\lambda_3$  is fully determined by  $\Lambda_1$  and  $\lambda_2$  as shown below.

$$\begin{aligned} \lambda_3 &= -\frac{1}{2}\mu^T\Sigma^{-1}\mu \\ &= -\frac{1}{2}(\Sigma\Sigma^{-1}\mu)^T(\Sigma^{-1}\mu) \\ &= -\frac{1}{2}\left((-2\Lambda_1)^{-1}\lambda_2\right)^T\lambda_2 \end{aligned} \quad (91)$$

We know that  $\mathcal{N}(\mathbf{z}|\mu, w\Sigma)\text{InvGam}(w|a, a)$  is a curved exponential family distribution since  $\lambda_3$  is fully determined by  $\Lambda_1$  and  $\lambda_2$ . ■

We give the update for the following model where the prior over  $p(\mathbf{z}, w)$  is a joint distribution of a Multivariate t-Distribution. Our method can be easily applied to models even when this is not the case. We use  $q(\mathbf{z}, w) = \mathcal{N}(\mathbf{z}|\mu, w\Sigma)\text{InvGam}(w|a, a)$  as the variational distribution. Note that when  $a \leq 1$ , the variance of the marginal distribution  $q(\mathbf{z})$  does not exist. We assume that  $a > 1$  such that the variance exists.

Given  $\mathbf{z} \in \mathcal{R}^d$  and  $w \in \mathcal{R}_{++}$ , we consider the following model :

$$p(\mathcal{D}, \mathbf{z}, w) = \prod_{n=1}^N p(\mathcal{D}_n|\mathbf{z}) \mathcal{N}(\mathbf{z}|\mathbf{0}, w\mathbf{I}) \text{InvGam}(w|a_0, a_0) \quad (92)$$

For this model, by using the Gauss-Newton approximation of the Hessian, the NGVI update given in Theorem 1 can be written as a weight-perturbed Newton method similar to Khan et al. (2018).

1. Generate a sample  $(\mathbf{z}_*, w_*) \sim q(\mathbf{z}, w|\lambda)$ .
2. Sample a random data example  $p(\mathcal{D}_n|\mathbf{z})$  and compute its gradient at  $\mathbf{z}_*$ .

$$\mathbf{g} \leftarrow \nabla_{\mathbf{z}} \log p(\mathcal{D}_n|\mathbf{z}) \quad (93)$$

3. Compute a *learning-rate* modifier:

$$u \leftarrow \frac{a}{(a + d/2 - 1)} \left[ 1 + \frac{1}{2} (\mathbf{z}_* - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{z}_* - \boldsymbol{\mu}) / a \right] \quad (94)$$

4. Update  $a$  for the Inverse-Gamma distribution:

$$a \leftarrow (1 - \beta_2)a + \beta_2 \left[ a_0 - \frac{Nw_*^2}{2(1 - w_*)} \mathbf{g}^\top \boldsymbol{\Sigma} \mathbf{g} \right], \quad (95)$$

where we choose  $\beta_2$  such that  $a > 1$ .

5. Take a Newton-like step to update  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ .

$$\mathbf{S} \leftarrow (1 - \beta_1)\mathbf{S} + \beta_1(u\mathbf{g}\mathbf{g}^\top) \quad (96)$$

$$\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} - \beta_1(\mathbf{S} + \mathbf{I}/N)^{-1}(\mathbf{g} + \boldsymbol{\mu}/N) \quad (97)$$

6. Repeat until convergence.

The update of  $\boldsymbol{\mu}$  and  $\mathbf{S}$  is very similar to the Gauss-Newton update of Khan et al. (2018), but here the learning rate  $\beta_1$  is multiplied by  $u$ . The algorithm therefore can be implemented easily by using further approximations discussed in Khan et al. (2018).

To derive this algorithm, we prove a version of Bonnet's and Price's theorem for the case when the expectation is taken with respect to a scale-mixture of Gaussian instead of a Gaussian.

First of all, we show how to update  $q(\mathbf{z}|w) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})$ . The natural parameter and the sufficient statistics of  $q(\mathbf{z}|w)$  are  $\boldsymbol{\lambda}_z = \{\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}, -\frac{1}{2}\boldsymbol{\Sigma}^{-1}\}$  and  $\phi_z(\mathbf{z}, w) = \{w^{-1}\mathbf{z}, w^{-1}\mathbf{z}\mathbf{z}^\top\}$ . The expectation parameters are

$$\mathbf{m}_{z_1} = \mathbb{E}_{q(z,w)} [w^{-1}\mathbf{z}] = \boldsymbol{\mu} \quad (98)$$

$$\mathbf{m}_{z_2} = \mathbb{E}_{q(z,w)} [w^{-1}\mathbf{z}\mathbf{z}^\top] = \boldsymbol{\mu}\boldsymbol{\mu}^\top + \boldsymbol{\Sigma} \quad (99)$$

By the chain rule, we know that

$$\nabla_{m_{z_1}} \mathcal{L}(\boldsymbol{\lambda}) = \nabla_{\boldsymbol{\mu}} \mathcal{L}(\boldsymbol{\lambda}) - 2\nabla_{\boldsymbol{\Sigma}} \mathcal{L}(\boldsymbol{\lambda})\boldsymbol{\mu} \quad (100)$$

$$\nabla_{m_{z_2}} \mathcal{L}(\boldsymbol{\lambda}) = \nabla_{\boldsymbol{\Sigma}} \mathcal{L}(\boldsymbol{\lambda}) \quad (101)$$

Let's denote  $h_n(\mathbf{z}) := -\log p(\mathcal{D}_n|\mathbf{z})$ . In the model, the lower bound  $\mathcal{L}(\boldsymbol{\lambda})$  can be expressed as

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(z,w)} \left[ -\sum_{n=1}^N h_n(\mathbf{z}) + \log \frac{\mathcal{N}(\mathbf{z}|\mathbf{0}, w\mathbf{I})}{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})} + \log \frac{\text{InvGam}(w|a_0, a_0)}{\text{InvGam}(w|a, a)} \right], \quad (102)$$

The gradients of the lower bound are

$$\nabla_{\boldsymbol{\mu}} \mathcal{L}(\boldsymbol{\lambda}) = \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(z,w)} \left[ -\sum_{n=1}^N h_n(\mathbf{z}) \right] - \boldsymbol{\mu} \quad (103)$$

$$\nabla_{\boldsymbol{\Sigma}} \mathcal{L}(\boldsymbol{\lambda}) = \nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{q(z,w)} \left[ -\sum_{n=1}^N h_n(\mathbf{z}) \right] - \frac{1}{2}\mathbf{I} + \frac{1}{2}\boldsymbol{\Sigma}^{-1} \quad (104)$$

The natural gradient update is

$$\boldsymbol{\lambda}_{z_1} \leftarrow \boldsymbol{\lambda}_{z_1} + \beta_1 \nabla_{m_{z_1}} \mathcal{L}(\boldsymbol{\lambda}) \quad (105)$$

$$\boldsymbol{\lambda}_{z_2} \leftarrow \boldsymbol{\lambda}_{z_2} + \beta_1 \nabla_{m_{z_2}} \mathcal{L}(\boldsymbol{\lambda}) \quad (106)$$

Since  $\boldsymbol{\lambda}_{z_1} = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$ ,  $\boldsymbol{\lambda}_{z_2} = -\frac{1}{2} \boldsymbol{\Sigma}^{-1}$ , using Eq.100,101,103,104, the update can be re-expressed as

$$\boldsymbol{\Sigma}^{-1} \leftarrow (1 - \beta_1) \boldsymbol{\Sigma}^{-1} + 2\beta_1 \nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{q(z,w)} \left[ \sum_{n=1}^N h_n(\mathbf{z}) \right] + \beta_1 \mathbf{I} \quad (107)$$

$$\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} - \beta_1 \boldsymbol{\Sigma} \left( \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(z,w)} \left[ \sum_{n=1}^N h_n(\mathbf{z}) \right] + \boldsymbol{\mu} \right) \quad (108)$$

By defining  $\mathbf{S} = [\boldsymbol{\Sigma}^{-1} - \mathbf{I}] / N$ , we have

$$\mathbf{S} \leftarrow (1 - \beta_1) \mathbf{S} + 2\beta_1 \nabla_{\mathbf{S}} \mathbb{E}_{q(z,w)} \left[ \sum_{n=1}^N h_n(\mathbf{z}) / N \right] \quad (109)$$

$$\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} - \beta_1 (\mathbf{S} + \mathbf{I}/N)^{-1} \left( \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(z,w)} \left[ \sum_{n=1}^N h_n(\mathbf{z}) / N \right] + \boldsymbol{\mu} / N \right) \quad (110)$$

Now, we show the update for  $q(w) = \text{InvGam}(w|a, a)$ . The natural parameter and the sufficient statistics of  $q(w)$  are  $\lambda_w = a$  and  $\phi_w(w) = -w^{-1} - \log w$ .

The gradient can be expressed as

$$\nabla_{m_w} \mathcal{L}(\boldsymbol{\lambda}) = \nabla_{m_w} \mathbb{E}_{q(z,w)} \left[ - \sum_{n=1}^N h_n(\mathbf{z}) \right] + a_0 - a \quad (111)$$

The update can be expressed in terms of  $a$  as

$$a \leftarrow (1 - \beta_2) a + \beta_2 \left( a_0 - \sum_{n=1}^N \nabla_{m_w} \mathbb{E}_{q(z,w)} [h_n(\mathbf{z})] \right) \quad (112)$$

While the gradient w.r.t. the expectation parameter does not admit a close-form expression, we can compute the gradient using the re-parametrization trick. Recall that the gradient  $\nabla_{m_w} \mathbb{E}_{q(z,w)} [h_n(\mathbf{z})]$  can be computed as

$$\nabla_{m_w} \mathbb{E}_{q(z,w)} [h(\mathbf{z})] = (\nabla_{\lambda_w} m_w)^{-1} \nabla_{\lambda_w} \mathbb{E}_{q(z,w)} [h_n(\mathbf{z})] \quad (113)$$

$$= (\nabla_{\lambda_w} \mathbb{E}_{q(w)} [\phi_w(w)])^{-1} \nabla_{\lambda_w} \mathbb{E}_{q(z,w)} [h_n(\mathbf{z})] \quad (114)$$

$$= (\nabla_a \mathbb{E}_{q(w)} [\phi_w(w)])^{-1} \nabla_a \mathbb{E}_{q(z,w)} [h_n(\mathbf{z})] \quad (115)$$

Note that  $\nabla_a \mathbb{E}_{q(w)} [\phi_w(w)] = \nabla_a m_w$  has a closed-form expression.

Even when  $\nabla_a m_w$  does not admit a closed-form expression, we can also compute it when  $q(w)$  is re-parameterizable shown below. Since  $q(w)$  is (implicitly) re-parameterizable, the gradient can be computed as

$$\nabla_a \mathbb{E}_{q(w)} [\phi_w(w)] = - \int \text{InvGam}(w|a, a) (\nabla_w [w^{-1} + \log w]) (\nabla_a w) dw \quad (116)$$

$$\nabla_a \mathbb{E}_{q(z,w)} [h_n(\mathbf{z})] = \int \int q(w) (\nabla_w q(\mathbf{z}|w)) (\nabla_a w) h_n(\mathbf{z}) dw d\mathbf{z} \quad (117)$$

$$= \int \int \text{InvGam}(w|a, a) (\nabla_w \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})) (\nabla_a w) h_n(\mathbf{z}) dw d\mathbf{z} \quad (118)$$

Now, we discuss about gradient approximation. Recall that the joint distribution is  $\text{InvGam}(w|a, a)\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})$ . The marginal distribution is a multivariate t-distribution as shown below

$$q(\mathbf{z}) = \det(\pi\boldsymbol{\Sigma})^{-1/2} \frac{\Gamma(a + d/2) \left(2a + (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu})\right)^{-a-d/2}}{\Gamma(a) (2a)^{-a}}. \quad (119)$$

When  $a > 1$ , we denote the following function as  $u(\mathbf{z})$

$$u(\mathbf{z}) := \int w q(\mathbf{z}|w) q(w) dw = \frac{a}{(a + d/2 - 1)} \left(1 + (2a)^{-1} (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu})\right) q(\mathbf{z}). \quad (120)$$

**Lemma 8** *Let  $h(\mathbf{z})$  be a twice continuously differentiable function and  $q(w) = \text{InvGam}(w|a, a)$ , where  $a > 1$ . If  $|h(\mathbf{z})|$ ,  $|z_j h(\mathbf{z})|$  and  $|\nabla_{z_j} h(\mathbf{z})|$  are integrable for any index  $j$  and  $k$ , the following gradient identities hold.*

$$\nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(w)\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})} [h(\mathbf{z})] = \int q(\mathbf{z}) \nabla_{\boldsymbol{\mu}} h(\mathbf{z}) d\mathbf{z} \quad (121)$$

$$\nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{q(w)\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})} [h(\mathbf{z})] = \frac{1}{2} \int u(\mathbf{z}) \nabla_{\boldsymbol{\Sigma}}^2 h(\mathbf{z}) d\mathbf{z} \quad (122)$$

The lemma can be generalized to the case when  $q(w)$  is a generalized inverse Gaussian distribution.

### Proof

$$\nabla_{\boldsymbol{\mu}_j} \mathbb{E}_{q(w)\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})} [h(\mathbf{z})] = \int \int q(w) [\nabla_{\boldsymbol{\mu}_j} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})] h(\mathbf{z}) dw d\mathbf{z} \quad (123)$$

$$= - \int q(w) \left\{ \int [\nabla_{z_j} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})] h(\mathbf{z}) d\mathbf{z} \right\} dw \quad (124)$$

$$= - \int q(w) \left\{ \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma}) h(\mathbf{z}) dz_{-j} \Big|_{z_j=-\infty}^{z_j=+\infty} - \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma}) \nabla_{z_j} h(\mathbf{z}) dz \right\} dw \quad (125)$$

$$= - \underbrace{\int q(z_j, \mathbf{z}_{-j}) h(z_j, \mathbf{z}_{-j}) dz_{-j} \Big|_{z_j=-\infty}^{z_j=+\infty}}_0 + \int q(w) \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma}) \nabla_{z_j} h(\mathbf{z}) dz dw \quad (126)$$

$$= \int q(\mathbf{z}) \nabla_{z_j} h(\mathbf{z}) d\mathbf{z} \quad (127)$$



We use integration by part to move from Eq. 124 to Eq. 125. Since  $|h(\mathbf{z})|$  is integrable, we know that the first term in Eq. 126 is 0.

Similarly, we have the following result by using integration by part twice.

$$\nabla_{\Sigma_{j,k}} \mathbb{E}_{q(w)\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\Sigma)} [h(\mathbf{z})] = \int q(w) [\nabla_{\Sigma_{j,k}} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\Sigma)] h(\mathbf{z}) dw d\mathbf{z} \quad (128)$$

$$= \frac{1}{2} \int \int q(w) [w \nabla_{z_j} \nabla_{z_k} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\Sigma)] h(\mathbf{z}) dw d\mathbf{z} \quad (129)$$

$$= \frac{1}{2} \underbrace{\int q(w) w \int [\nabla_{z_k} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\Sigma)] h(\mathbf{z}) d\mathbf{z}_{-j} dw}_{0} \Big|_{z_j=-\infty}^{z_j=+\infty} - \frac{1}{2} \int \int q(w) w \nabla_{z_k} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\Sigma) \nabla_{z_j} h(\mathbf{z}) d\mathbf{z} dw \quad (130)$$

$$= -\frac{1}{2} \int q(w) w \int \nabla_{z_k} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\Sigma) \nabla_{z_j} h(\mathbf{z}) d\mathbf{z} dw \quad (131)$$

$$= -\frac{1}{2} \int q(w) w \left\{ \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\Sigma) \nabla_{z_j} h(\mathbf{z}) d\mathbf{z}_{-k} \Big|_{z_k=-\infty}^{z_k=+\infty} - \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\Sigma) \nabla_{z_k} \nabla_{z_j} h(\mathbf{z}) d\mathbf{z} \right\} dw \quad (132)$$

$$= -\frac{1}{2} \underbrace{\int u(z_k, \mathbf{z}_{-k}) \nabla_{z_j} h(\mathbf{z}) d\mathbf{z}_{-k} \Big|_{z_k=-\infty}^{z_k=+\infty}}_0 + \frac{1}{2} \int q(w) w \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\Sigma) \nabla_{z_k} \nabla_{z_j} h(\mathbf{z}) d\mathbf{z} dw \quad (133)$$

$$= \frac{1}{2} \int u(\mathbf{z}) \nabla_{z_k} \nabla_{z_j} h(\mathbf{z}) d\mathbf{z} \quad (134)$$

We use integration by part twice to move from Eq. 128 to Eq. 130 and from Eq. 131 to Eq. 132. The first term in Eq. 130 is 0 since  $|z_k h(\mathbf{z})|$  is integrable. The first term in Eq. 133 is 0 since  $|z_k h(\mathbf{z})|$  and  $|\nabla_{z_j} h(\mathbf{z})|$  are integrable.  $\blacksquare$

Using Lemma 8, the update for  $q(\mathbf{z}|w)$  can be expressed as

$$\mathbf{S} \leftarrow (1 - \beta_1) \mathbf{S} + \beta_1 \mathbb{E}_{q(z)} \left[ \frac{a}{(a + d/2 - 1)} \left( 1 + (2a)^{-1} (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu}) \right) \left( \sum_{n=1}^N \nabla_z^2 h_n(\mathbf{z}) / N \right) \right] \quad (135)$$

$$\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} - \beta_1 (\mathbf{S} + \mathbf{I}/N)^{-1} \left( \mathbb{E}_{q(z)} \left[ \sum_{n=1}^N \nabla_z h_n(\mathbf{z}) / N \right] + \boldsymbol{\mu} / N \right) \quad (136)$$

Using MC approximation, we have the following update

$$\mathbf{S} \leftarrow (1 - \beta_1) \mathbf{S} + \beta_1 \frac{a}{(a + d/2 - 1)} \left( 1 + (2a)^{-1} (\mathbf{z}_* - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z}_* - \boldsymbol{\mu}) \right) \left( \sum_{n=1}^N \nabla_{z_*}^2 h_n(\mathbf{z}) / N \right) \quad (137)$$

$$\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} - \beta_1 (\mathbf{S} + \mathbf{I}/N)^{-1} \left( \sum_{n=1}^N \nabla_{z_*} h_n(\mathbf{z}) / N + \boldsymbol{\mu} / N \right) \quad (138)$$

Similarly, the update for  $q(w)$  can be expressed as

$$a \leftarrow (1 - \beta_2)a + \beta_2 \left( a_0 - \frac{N}{\nabla_a \mathbb{E}_{q(w)} [\phi_w(w)]} \sum_{n=1}^N \nabla_a \mathbb{E}_{q(z,w)} [h_n(\mathbf{z})] / N \right) \quad (139)$$

where

$$\nabla_a \mathbb{E}_{q(w)} [\phi_w(w)] = - \int \text{InvGam}(w|a, a) (\nabla_w [w^{-1} + \log w]) (\nabla_a w) dw \quad (140)$$

$$\nabla_a \mathbb{E}_{q(z,w)} [h_n(\mathbf{z})] = \int \int \text{InvGam}(w|a, a) (\nabla_w \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})) (\nabla_a w) h_n(\mathbf{z}) dw d\mathbf{z} \quad (141)$$

Using one MC sample, we have

$$\nabla_a \mathbb{E}_{q(w)} [\phi_w(w)] \approx (w_*^{-2} - w_*^{-1}) \nabla_a w_* \quad (142)$$

$$\nabla_a \mathbb{E}_{q(z,w)} [h_n(\mathbf{z})] \approx \int (\nabla_{w_*} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w_*\boldsymbol{\Sigma})) (\nabla_a w_*) h_n(\mathbf{z}) d\mathbf{z} \quad (143)$$

$$= \int \text{Tr} \left( \boldsymbol{\Sigma} \nabla_{\hat{\boldsymbol{\Sigma}}} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \hat{\boldsymbol{\Sigma}}) \right) (\nabla_a w_*) h_n(\mathbf{z}) d\mathbf{z} \quad (144)$$

$$= \nabla_a w_* \text{Tr} \left( \boldsymbol{\Sigma} \nabla_{\hat{\boldsymbol{\Sigma}}} \mathbb{E}_{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \hat{\boldsymbol{\Sigma}})} [h_n(\mathbf{z})] \right) \quad (145)$$

$$= \frac{\nabla_a w_*}{2} \text{Tr} \left( \boldsymbol{\Sigma} \mathbb{E}_{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \hat{\boldsymbol{\Sigma}})} [\nabla_{\mathbf{z}}^2 h_n(\mathbf{z})] \right) \quad (146)$$

$$\approx \frac{\nabla_a w_*}{2} \text{Tr} \left( \boldsymbol{\Sigma} \nabla_{z_*}^2 h_n(\mathbf{z}) \right), \quad (147)$$

where  $\hat{\boldsymbol{\Sigma}} = w_* \boldsymbol{\Sigma}$  and we use the Price's theorem  $\nabla_{\hat{\boldsymbol{\Sigma}}} \mathbb{E}_{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \hat{\boldsymbol{\Sigma}})} [h_n(\mathbf{z})] = \frac{1}{2} \mathbb{E}_{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \hat{\boldsymbol{\Sigma}})} [\nabla_{\mathbf{z}}^2 h_n(\mathbf{z})]$ .

## Appendix D. Updates for Matrix-Variate Gaussian Distribution

We first show that MVG is a multi-linear exponential-family distribution.

**Lemma 9** *Matrix Gaussian distribution is a member of the multi-linear exponential family*

**Proof** Let  $\boldsymbol{\Lambda}_1 = \mathbf{W}$ ,  $\boldsymbol{\Lambda}_2 = \mathbf{U}^{-1}$ , and  $\boldsymbol{\Lambda}_3 = \mathbf{V}^{-1}$ . The distribution on  $\mathbf{Z} \in \mathcal{R}^{d \times p}$  can be expressed as follows.

$$\mathcal{MN}(\mathbf{Z}|\mathbf{W}, \mathbf{U}, \mathbf{V}) \quad (148)$$

$$= (2\pi)^{-dp/2} \exp \left[ -\frac{1}{2} \text{Tr} \left( \mathbf{V}^{-1} (\mathbf{Z} - \mathbf{W})^T \mathbf{U}^{-1} (\mathbf{Z} - \mathbf{W}) \right) - (d/2 \log \text{Det}(\mathbf{V}) + p/2 \log \text{Det}(\mathbf{U})) \right]$$

$$= (2\pi)^{-dp/2} \exp \left\{ \text{Tr} \left( \boldsymbol{\Lambda}_3 \left( -\frac{1}{2} \mathbf{Z} + \boldsymbol{\Lambda}_1 \right)^T \boldsymbol{\Lambda}_2 \mathbf{Z} \right) - \frac{1}{2} \left[ \text{Tr} \left( \boldsymbol{\Lambda}_3 \boldsymbol{\Lambda}_1^T \boldsymbol{\Lambda}_2 \boldsymbol{\Lambda}_1 \right) + d \log \text{Det}(\boldsymbol{\Lambda}_3) + p \log \text{Det}(\boldsymbol{\Lambda}_2) \right] \right\}. \quad (149)$$

The function  $\text{Tr} \left( \boldsymbol{\Lambda}_3 \left( -\frac{1}{2} \mathbf{Z} + \boldsymbol{\Lambda}_1 \right)^T \boldsymbol{\Lambda}_2 \mathbf{Z} \right)$  is linear with respect each  $\boldsymbol{\Lambda}_j$  given others.  $\blacksquare$

We now derive the NGVI update using our new expectation parameterization. We can obtain function  $\phi_1$ ,  $\phi_2$ , and  $\phi_3$  from the multi-linear function

$$f(\mathbf{Z}, \mathbf{\Lambda}) := \text{Tr} \left( \mathbf{\Lambda}_3 \left( -\frac{1}{2} \mathbf{Z} + \mathbf{\Lambda}_1 \right)^T \mathbf{\Lambda}_2 \mathbf{Z} \right). \quad (150)$$

For example, we can obtain function  $\phi_1$  from  $f(\mathbf{Z}, \mathbf{\Lambda})$  as shown below:

$$f(\mathbf{Z}, \mathbf{\Lambda}) = \underbrace{\langle \mathbf{\Lambda}_1, \mathbf{\Lambda}_2 \mathbf{Z} \mathbf{\Lambda}_3 \rangle}_{\phi_1(\mathbf{Z}, \mathbf{\Lambda}_{-1})} - \underbrace{\frac{1}{2} \text{Tr} \left( \mathbf{\Lambda}_3 \mathbf{Z}^T \mathbf{\Lambda}_2 \mathbf{Z} \right)}_{r_1(\mathbf{Z}, \mathbf{\Lambda}_{-1})}. \quad (151)$$

Similarly, we can obtain functions  $\phi_2$  and  $\phi_3$ . The corresponding expectation parameters of the Matrix Gaussian distribution can then be derived as below:

$$\mathbf{M}_1 = \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [\mathbf{\Lambda}_2 \mathbf{Z} \mathbf{\Lambda}_3] = \mathbf{\Lambda}_2 \mathbf{\Lambda}_1 \mathbf{\Lambda}_3 \quad (152)$$

$$\mathbf{M}_2 = \mathbb{E}_{\mathcal{MN}(z|w,u,v)} \left[ -\frac{1}{2} \mathbf{Z} \mathbf{\Lambda}_3 \mathbf{Z}^T + \mathbf{Z} \mathbf{\Lambda}_3 \mathbf{\Lambda}_1^T \right] = \frac{1}{2} (\mathbf{\Lambda}_1 \mathbf{\Lambda}_3 \mathbf{\Lambda}_1^T - p \mathbf{\Lambda}_2^{-1}) \quad (153)$$

$$\mathbf{M}_3 = \mathbb{E}_{\mathcal{MN}(z|w,u,v)} \left[ -\frac{1}{2} \mathbf{Z}^T \mathbf{\Lambda}_2 \mathbf{Z} + \mathbf{\Lambda}_1^T \mathbf{\Lambda}_2 \mathbf{Z} \right] = \frac{1}{2} (\mathbf{\Lambda}_1^T \mathbf{\Lambda}_2 \mathbf{\Lambda}_1 - d \mathbf{\Lambda}_3^{-1}) \quad (154)$$

We can then compute the gradient with respect to the expectation parameters using chain-rule:

$$\nabla_{M_1} \mathbb{E}_{q(z|\lambda)} [h(\mathbf{Z})] = (\mathbf{\Lambda}_2)^{-1} \nabla_W \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [h(\mathbf{Z})] (\mathbf{\Lambda}_3)^{-1} \quad (155)$$

$$\nabla_{M_2} \mathbb{E}_{q(z|\lambda)} [h(\mathbf{Z})] = \frac{-2}{p} \nabla_U \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [h(\mathbf{Z})] \quad (156)$$

$$\nabla_{M_3} \mathbb{E}_{q(z|\lambda)} [h(\mathbf{Z})] = \frac{-2}{d} \nabla_V \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [h(\mathbf{Z})] \quad (157)$$

We will now express the gradients in terms of the gradient of the function  $h(\mathbf{Z})$ . This leads to a simple update because gradient of  $h(\mathbf{Z})$  can be obtained using automatic gradients (or backpropagation when using a neural network). Let  $\mathbf{z} = \text{vec}(\mathbf{Z})$  and  $\mathbf{Z} = \text{Mat}(\mathbf{z})$ . The distribution can be re-expressed as a multivariate Gaussian distribution  $\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\mu} = \text{vec}(\mathbf{W})$ ,  $\boldsymbol{\Sigma} = \mathbf{V} \otimes \mathbf{U}$ , and  $\otimes$  denotes the Kronecker product. Furthermore, the lower bound can be re-expressed as  $\mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[ -\hat{h}(\mathbf{z}) \right]$ , where  $\hat{h}(\mathbf{z}) = h(\mathbf{Z})$ . We make use of the Bonnet's and Price's theorems (Opper and Archambeau, 2009):

$$\nabla_{\boldsymbol{\mu}} \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[ \hat{h}(\mathbf{z}) \right] = \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[ \nabla_z \hat{h}(\mathbf{z}) \right] \quad (158)$$

$$\nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[ \nabla_z \hat{h}(\mathbf{z}) \right] = \frac{1}{2} \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[ \nabla_z^2 \hat{h}(\mathbf{z}) \right] \quad (159)$$

These identities can be used to express the gradient with respect to the expectation parameters in terms of the gradient with respect to  $\mathbf{Z}$ :

$$\nabla_W \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [h(\mathbf{Z})] = \text{Mat} \left( \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[ \nabla_z \hat{h}(\mathbf{z}) \right] \right) \quad (160)$$

$$= \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [\nabla_Z h(\mathbf{Z})] \quad (161)$$

$$\nabla_U \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [h(\mathbf{Z})] = (\nabla_U \Sigma) \nabla_\Sigma \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[ \nabla_z \hat{h}(\mathbf{z}) \right] \quad (162)$$

$$= \frac{1}{2} (\nabla_U \Sigma) \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[ \nabla_z^2 \hat{h}(\mathbf{z}) \right] \quad (163)$$

$$\approx \frac{1}{2} (\nabla_U \Sigma) \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[ \nabla_z \hat{h}(\mathbf{z}) \nabla_z \hat{h}(\mathbf{z})^T \right] \quad (164)$$

$$= \frac{1}{2} \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [\nabla_Z h(\mathbf{Z}) \mathbf{V} \nabla_Z h(\mathbf{Z})^T] \quad (165)$$

$$\nabla_V \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [h(\mathbf{Z})] = (\nabla_V \Sigma) \nabla_\Sigma \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[ \nabla_z \hat{h}(\mathbf{z}) \right] \quad (166)$$

$$= \frac{1}{2} (\nabla_V \Sigma) \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[ \nabla_z^2 \hat{h}(\mathbf{z}) \right] \quad (167)$$

$$\approx \frac{1}{2} (\nabla_V \Sigma) \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[ \nabla_z \hat{h}(\mathbf{z}) \nabla_z \hat{h}(\mathbf{z})^T \right] \quad (168)$$

$$= \frac{1}{2} \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [\nabla_Z h(\mathbf{Z})^T \mathbf{U} \nabla_Z h(\mathbf{Z})]. \quad (169)$$

To avoid computation of the Hessian, we have used the Gauss-Newton approximation in Eq 164 and Eq. 168 (Khan et al., 2018).

We choose the step-size as  $\beta = \{\beta_1, p\beta_2, d\beta_2\}$ . The update with the Gauss-Newton approximation can be expressed as

$$\mathbf{\Lambda}_1 \leftarrow \mathbf{\Lambda}_1 - \beta_1 (\mathbf{\Lambda}_2)^{-1} \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [\nabla_Z h(\mathbf{Z})] (\mathbf{\Lambda}_3)^{-1} \quad (170)$$

$$\mathbf{\Lambda}_2 \leftarrow \mathbf{\Lambda}_2 + \beta_2 \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [\nabla_Z h(\mathbf{Z}) \mathbf{V} \nabla_Z h(\mathbf{Z})^T] \quad (171)$$

$$\mathbf{\Lambda}_3 \leftarrow \mathbf{\Lambda}_3 + \beta_2 \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [\nabla_Z h(\mathbf{Z})^T \mathbf{U} \nabla_Z h(\mathbf{Z})] \quad (172)$$

We can re-express these in terms of  $\{\mathbf{W}, \mathbf{U}^{-1}, \mathbf{V}^{-1}\}$  to get the final updates:

$$\mathbf{W} \leftarrow \mathbf{W} - \beta_1 \mathbf{U} \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [\nabla_Z h(\mathbf{Z})] \mathbf{V} \quad (173)$$

$$(\mathbf{U})^{-1} \leftarrow (\mathbf{U})^{-1} + \beta_2 \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [\nabla_Z h(\mathbf{Z}) \mathbf{V} \nabla_Z h(\mathbf{Z})^T] \quad (174)$$

$$(\mathbf{V})^{-1} \leftarrow (\mathbf{V})^{-1} + \beta_2 \mathbb{E}_{\mathcal{MN}(z|w,u,v)} [\nabla_Z h(\mathbf{Z})^T \mathbf{U} \nabla_Z h(\mathbf{Z})] \quad (175)$$

## Appendix E. Connection to Block Mirror Descent

As discussed in Section 4, the following theorem shows that the block natural-gradient update is indeed a block mirror descent.

**Theorem 10** *The block natural-gradient update is a block mirror descent.*

In fact, we can perform block natural-gradient update for the class of mixture of exponential family. The following lemmas give the proof-sketch of the theorem. Recall that  $j$  is the block index of a variational distribution studied in this paper. We first define the function  $\Psi_j(\boldsymbol{\lambda}_j) = \mathbb{E}_{q(\cdot|\lambda_j, \lambda_{-j})} [A_j(\cdot)]$ , where  $A_j(\cdot)$  is the log-partition function of the distribution at

block  $j$ . As shown in the following examples,  $\Psi_j(\boldsymbol{\lambda}_j)$  is convex w.r.t.  $\boldsymbol{\lambda}_j$  given  $\boldsymbol{\lambda}_{-j}$  is known. If  $\Psi_j(\boldsymbol{\lambda}_j)$  is strictly convex w.r.t.  $\boldsymbol{\lambda}_j$ , we say the minimality holds. We then define the convex conjugate of function  $\Psi_j$  as below.

$$\Psi_j^*(\mathbf{x}) := \sup_{\boldsymbol{\eta}} v_j(\mathbf{x}, \boldsymbol{\eta}), \quad (176)$$

where  $v_j(\boldsymbol{\eta}) = \langle \mathbf{x}, \boldsymbol{\eta} \rangle - \Psi_j(\boldsymbol{\eta})$ .

**Example 4 (Mixture of exponential family distribution)** *The variational distribution is a mixture of exponential family distribution:  $q(\mathbf{z}, \mathbf{w} | \boldsymbol{\lambda}_z, \boldsymbol{\lambda}_w)$ . When  $j = w$ , the log-partition function is  $A_j(\boldsymbol{\lambda}_j) = A_w(\boldsymbol{\lambda}_w)$ . The expectation parameter is  $\mathbf{m}_j = \mathbb{E}_q[\boldsymbol{\phi}_w(\mathbf{w})]$ . Since  $q(\mathbf{w})$  is an exponential family distribution when  $j = w$ , we know that  $A_w(\boldsymbol{\lambda}_w)$  is convex w.r.t.  $\boldsymbol{\lambda}_w$ . Furthermore,  $\Psi_j(\boldsymbol{\lambda}_j) = A_w(\boldsymbol{\lambda}_w)$  is convex w.r.t.  $\boldsymbol{\lambda}_j$  given  $\boldsymbol{\lambda}_{-j}$  is known. The minimal condition implies that  $q(\mathbf{w})$  is a minimal exponential family distribution. When  $j = z$ , the log-partition function is  $A_j(\boldsymbol{\lambda}_j) = A_z(\boldsymbol{\lambda}_z, \mathbf{w})$ . The expectation parameter is  $\mathbf{m}_j = \mathbb{E}_q[\boldsymbol{\phi}_z(\mathbf{z}, \mathbf{w})]$ . According Eq. 47, we know that  $A_z(\boldsymbol{\lambda}_z, \mathbf{w})$  is convex w.r.t.  $\boldsymbol{\lambda}_j$  for any valid  $\mathbf{w}$ , which implies that  $\Psi_j(\boldsymbol{\lambda}_j) = \mathbb{E}_{q(w)}[A_z(\boldsymbol{\lambda}_z, \mathbf{w})]$  is convex w.r.t.  $\boldsymbol{\lambda}_j$  given  $\boldsymbol{\lambda}_{-j}$  is known. If Lemma 5 holds, we know that the minimal condition holds.*

**Example 5 (Multi-linear exponential family distribution)** *The variational distribution is a multi-linear exponential family distribution:  $q(\mathbf{z} | \boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N)$ . When  $j \in \{1, \dots, N\}$ , the log-partition function is  $A_j(\boldsymbol{\lambda}_j) = A_z(\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j})$ . The expectation parameter is  $\mathbf{m}_j = \mathbb{E}_q[\boldsymbol{\phi}_j(\mathbf{z}, \boldsymbol{\lambda}_{-j})]$ . Since  $\boldsymbol{\lambda}_{-j}$  is known, we know that  $q(\mathbf{z} | \boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}) = q_j(\mathbf{z} | \boldsymbol{\lambda}_j)$  becomes a one-parameter exponential family distribution. We know that  $\Psi_j(\boldsymbol{\lambda}_j) = A_z(\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j})$  is convex w.r.t.  $\boldsymbol{\lambda}_j$  given that  $\boldsymbol{\lambda}_{-j}$  is known. The minimal condition implies that for all  $j \in \{1, \dots, N\}$ ,  $q_j(\mathbf{z} | \boldsymbol{\lambda}_j)$  is a minimal exponential family distribution.*

**Lemma 11** *If the minimality holds, we have*

$$\Psi_j^*(\mathbf{m}_j) = \sup_{\boldsymbol{\eta}} v_j(\mathbf{m}_j, \boldsymbol{\eta}) \quad (177)$$

$$= \langle \mathbf{m}_j, \boldsymbol{\lambda}_j \rangle - \Psi_j(\boldsymbol{\lambda}_j), \quad (178)$$

where  $\boldsymbol{\lambda}_j$  and  $\mathbf{m}_j$  are the natural parameter and the expectation parameter of the variational distribution at block  $j$ .

**Proof** We give the proof for the case when the variational distribution is a mixture of exponential family distribution. It is easy to show the lemma holds in the multi-linear exponential family case.

When  $j = w$ , since  $q(\mathbf{w})$  is an exponential family distribution, we know that

$$\mathbf{m}_w = \nabla_{\boldsymbol{\lambda}_w} A_w(\boldsymbol{\lambda}_w), \quad (179)$$

which implies that  $\mathbf{m}_w - \nabla_{\boldsymbol{\lambda}_w} A_w(\boldsymbol{\lambda}_w) = \mathbf{0}$ .

Note that the gradient of  $v_j(\mathbf{m}_w, \boldsymbol{\eta})$  is

$$\nabla_{\boldsymbol{\eta}} v_j(\mathbf{m}_w, \boldsymbol{\eta}) = \nabla_{\boldsymbol{\eta}} \langle \mathbf{m}_w, \boldsymbol{\eta} \rangle - \Psi_j(\boldsymbol{\eta}) = \mathbf{m}_w - \nabla_{\boldsymbol{\eta}} A_w(\boldsymbol{\eta}). \quad (180)$$

When  $\boldsymbol{\eta} = \boldsymbol{\lambda}_w$ , the gradient is  $\mathbf{0}$ . Since  $\Psi_j(\boldsymbol{\lambda}_w)$  is convex w.r.t.  $\boldsymbol{\lambda}_w$ , we know that  $\boldsymbol{\lambda}_w$  is an optimal solution. Since the minimality holds, we know that  $\Psi_j(\boldsymbol{\lambda}_w)$  is strictly convex, which implies that the solution is unique. Therefore,  $\underbrace{\Psi_w^*(\mathbf{m}_w)}_{\Psi_j^*(\mathbf{m}_j)} = \underbrace{\langle \mathbf{m}_w, \boldsymbol{\lambda}_w \rangle - A_w(\boldsymbol{\lambda}_w)}_{\langle \mathbf{m}_j, \boldsymbol{\lambda}_j \rangle - \Psi_j(\boldsymbol{\lambda}_j)}$ .

When  $j = z$ , by Eq. 21, we know that

$$\mathbf{m}_z = \nabla_{\lambda_z} \mathbb{E}_{q(w)} [A_z(\boldsymbol{\lambda}_z, \mathbf{w})], \quad (181)$$

which implies that  $\mathbf{m}_z - \nabla_{\lambda_z} \mathbb{E}_{q(w)} [A_z(\boldsymbol{\lambda}_z, \mathbf{w})] = \mathbf{0}$ . Similarly, since the minimality holds, when  $\boldsymbol{\lambda}_w$  is known, we know  $\underbrace{\Psi_z^*(\mathbf{m}_z)}_{\Psi_j^*(\mathbf{m}_j)} = \underbrace{\langle \mathbf{m}_z, \boldsymbol{\lambda}_z \rangle - \mathbb{E}_q [A_z(\boldsymbol{\lambda}_z, \mathbf{w})]}_{\langle \mathbf{m}_j, \boldsymbol{\lambda}_j \rangle - \Psi_j(\boldsymbol{\lambda}_j)}$ .  $\blacksquare$

Now, we show that the connection between block mirror descent update and the natural-gradient update.

**Lemma 12** *If the minimality holds, given that  $\boldsymbol{\lambda}_{-j} = \boldsymbol{\lambda}_{-j}^k$ , the following identity is true.*

$$\mathbb{B}_{\Psi_j^*}(\mathbf{m}_j \| \mathbf{m}_j^k) = \mathbb{D}_{KL}(q(\cdot | \boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k) \| q(\cdot | \boldsymbol{\lambda}_j^k, \boldsymbol{\lambda}_{-j}^k)) \quad (182)$$

$$= \Psi_j(\boldsymbol{\lambda}_j^k) - \Psi_j(\boldsymbol{\lambda}_j) - \left\langle \nabla_{\lambda_j} \Psi_j(\boldsymbol{\lambda}_j), \left( \boldsymbol{\lambda}_j^k - \boldsymbol{\lambda}_j \right) \right\rangle, \quad (183)$$

where  $\mathbb{B}(\cdot)$  denotes the Bregman divergence.

### Proof

We give the proof for the case when the variational distribution is a mixture of exponential family distribution. It is easy to show the lemma holds in the multi-linear exponential family case. According to the duality of the Bregman divergence and Lemma 11, we have

$$\mathbb{B}_{\Psi_j^*}(\mathbf{m}_j \| \mathbf{m}_j^k) = \mathbb{B}_{\Psi_j}(\boldsymbol{\lambda}_j^k \| \boldsymbol{\lambda}_j) \quad (184)$$

When  $j = w$ ,  $\Psi_w(\boldsymbol{\lambda}_w) = A_w(\boldsymbol{\lambda}_w)$ . Therefore, we have the following results.

$$\mathbb{D}_{KL}(q(\mathbf{z}, \mathbf{w} | \boldsymbol{\lambda}_w, \boldsymbol{\lambda}_z^k) \| q(\mathbf{z}, \mathbf{w} | \boldsymbol{\lambda}_w^k, \boldsymbol{\lambda}_z^k)) \quad (185)$$

$$= \mathbb{D}_{KL}(q(\mathbf{w} | \boldsymbol{\lambda}_w) \| q(\mathbf{w} | \boldsymbol{\lambda}_w^k)) \quad (186)$$

$$= \left\langle \mathbb{E}_{q(w|\lambda_w)} [\boldsymbol{\phi}_w(\mathbf{w})], \left( \boldsymbol{\lambda}_w - \boldsymbol{\lambda}_w^k \right) \right\rangle - A_w(\boldsymbol{\lambda}_w) + A_w(\boldsymbol{\lambda}_w^k) \quad (187)$$

$$= \left\langle \mathbf{m}_w, \left( \boldsymbol{\lambda}_w - \boldsymbol{\lambda}_w^k \right) \right\rangle - A_w(\boldsymbol{\lambda}_w) + A_w(\boldsymbol{\lambda}_w^k) \quad (188)$$

$$= \Psi_w(\boldsymbol{\lambda}_w^k) - \Psi_w(\boldsymbol{\lambda}_w) - \left\langle \nabla_{\lambda_w} \Psi_w(\boldsymbol{\lambda}_w), \left( \boldsymbol{\lambda}_w^k - \boldsymbol{\lambda}_w \right) \right\rangle \quad (189)$$

$$= \mathbb{B}_{\Psi_w}(\boldsymbol{\lambda}_w^k \| \boldsymbol{\lambda}_w), \quad (190)$$

where  $\mathbf{m}_w = \nabla_{\lambda_w} A_w(\boldsymbol{\lambda}_w) = \nabla_{\lambda_w} \Psi_w(\boldsymbol{\lambda}_w)$ .

When  $j = z$ ,  $\Psi_z(\boldsymbol{\lambda}_z) = \mathbb{E}_{q(w|\boldsymbol{\lambda}_z^k)} [A_z(\boldsymbol{\lambda}_z, \mathbf{w})]$ . Therefore, we have the following results.

$$\mathbb{D}_{KL}(q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda}_z, \boldsymbol{\lambda}_w^k) \| q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda}_z^k, \boldsymbol{\lambda}_w^k)) \quad (191)$$

$$= \mathbb{D}_{KL}(q(\mathbf{w}|\boldsymbol{\lambda}_w^k)q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) \| q(\mathbf{w}|\boldsymbol{\lambda}_w^k)q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z^k)) \quad (192)$$

$$= \left\langle \mathbb{E}_{q(z,w|\lambda_z, \lambda_w^k)} [\phi_z(\mathbf{z}, \mathbf{w})], \left( \boldsymbol{\lambda}_z - \boldsymbol{\lambda}_z^k \right) \right\rangle - \mathbb{E}_{q(w|\lambda_w^k)} [A_z(\boldsymbol{\lambda}_z, \mathbf{w})] + \mathbb{E}_{q(w|\lambda_w^k)} [A_z(\boldsymbol{\lambda}_z^k, \mathbf{w})] \quad (193)$$

$$= \left\langle \mathbf{m}_z, \left( \boldsymbol{\lambda}_z - \boldsymbol{\lambda}_z^k \right) \right\rangle - \Psi_z(\boldsymbol{\lambda}_z) + \Psi_z(\boldsymbol{\lambda}_z^k) \quad (194)$$

$$= \Psi_z(\boldsymbol{\lambda}_z^k) - \Psi_z(\boldsymbol{\lambda}_z) - \left\langle \nabla_{\lambda_z} \Psi_z(\boldsymbol{\lambda}_z), \left( \boldsymbol{\lambda}_z^k - \boldsymbol{\lambda}_z \right) \right\rangle \quad (195)$$

$$= \mathbb{B}_{\Psi_z}(\boldsymbol{\lambda}_z^k \| \boldsymbol{\lambda}_z), \quad (196)$$

where  $\mathbf{m}_z = \nabla_{\lambda_z} \Psi_z(\boldsymbol{\lambda}_z)$  due to Eq. 181. ■

**Lemma 13** *If the minimality holds, given  $\boldsymbol{\lambda}_{-j} = \boldsymbol{\lambda}_{-j}^k$ , a valid step-size  $\beta_j^k$ , and the following optimization problem*

$$\max_{\mathbf{m}_j} \left\langle \mathbf{m}_j, \widehat{\nabla}_{m_j^k} \mathcal{L}(\boldsymbol{\lambda}) \right\rangle - \frac{1}{\beta_j^k} \mathbb{B}_{\Psi_j^*}(\mathbf{m}_j \| \mathbf{m}_j^k), \quad (197)$$

the solution of the optimization problem is

$$\boldsymbol{\lambda}_j^{k+1} = \boldsymbol{\lambda}_j^k + \beta_j^k \widehat{\nabla}_{m_j^k} \mathcal{L}(\boldsymbol{\lambda}) \quad (198)$$

**Proof** According to Lemma 12, Eq. 197 can be also re-written as

$$\left\langle \mathbf{m}_j, \widehat{\nabla}_{m_j^k} \mathcal{L}(\boldsymbol{\lambda}) \right\rangle - \frac{1}{\beta_j^k} \mathbb{B}_{\Psi_j^*}(\mathbf{m}_j \| \mathbf{m}_j^k) \quad (199)$$

$$= \left\langle \mathbf{m}_j, \widehat{\nabla}_{m_j^k} \mathcal{L}(\boldsymbol{\lambda}) \right\rangle - \frac{1}{\beta_j^k} \mathbb{D}_{KL}(q(\cdot|\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k) \| q(\cdot|\boldsymbol{\lambda}_j^k, \boldsymbol{\lambda}_{-j}^k)) \quad (200)$$

$$= \left\langle \mathbf{m}_j, \widehat{\nabla}_{m_j^k} \mathcal{L}(\boldsymbol{\lambda}) \right\rangle - \frac{1}{\beta_j^k} \left\{ \left\langle \mathbf{m}_j, \left( \boldsymbol{\lambda}_j - \boldsymbol{\lambda}_j^k \right) \right\rangle - \Psi_j(\boldsymbol{\lambda}_j) + \Psi_j(\boldsymbol{\lambda}_j^k) \right\} \quad (201)$$

Similarly, by taking the gradient of Eq. 201 w.r.t.  $\mathbf{m}_j$  to be  $\mathbf{0}$ , we have

$$\begin{aligned} \mathbf{0} &= \nabla_{m_j} \left\{ \left\langle \mathbf{m}_j, \widehat{\nabla}_{m_j^k} \mathcal{L}(\boldsymbol{\lambda}) \right\rangle - \frac{1}{\beta_j^k} \left\{ \left\langle \mathbf{m}_j, \left( \boldsymbol{\lambda}_j - \boldsymbol{\lambda}_j^k \right) \right\rangle - \Psi_j(\boldsymbol{\lambda}_j) + \Psi_j(\boldsymbol{\lambda}_j^k) \right\} \right\} \\ &= \widehat{\nabla}_{m_j^k} \mathcal{L}(\boldsymbol{\lambda}) - \frac{1}{\beta_j^k} \left[ \left( \boldsymbol{\lambda}_j - \boldsymbol{\lambda}_j^k \right) + \left( \nabla_{m_j} \boldsymbol{\lambda}_j \right) \mathbf{m}_j - \nabla_{m_j} \Psi_j(\boldsymbol{\lambda}_j) \right] \end{aligned} \quad (202)$$

$$= \widehat{\nabla}_{m_j^k} \mathcal{L}(\boldsymbol{\lambda}) - \frac{1}{\beta_j^k} \left[ \left( \boldsymbol{\lambda}_j - \boldsymbol{\lambda}_j^k \right) + \left( \nabla_{m_j} \boldsymbol{\lambda}_j \right) \mathbf{m}_j - \underbrace{\left( \nabla_{\lambda_j} \Psi_j(\boldsymbol{\lambda}_j) \right)}_{\mathbf{m}_j} \right] \quad (203)$$

$$= \widehat{\nabla}_{m_j^k} \mathcal{L}(\boldsymbol{\lambda}) - \frac{1}{\beta_j^k} \left( \boldsymbol{\lambda}_j - \boldsymbol{\lambda}_j^k \right) \quad (204)$$

Finally, we obtain the following solution by solving Eq. 204

$$\boldsymbol{\lambda}_j^{k+1} = \boldsymbol{\lambda}_j^k + \beta_j^k \widehat{\nabla}_{m_j} \mathcal{L}(\boldsymbol{\lambda}) \quad (205)$$

■

## Appendix F. A Variant of Block Mirror Descent

Now, we discuss a variant of block mirror descent. Instead of sequentially performing update block by block, we can simultaneously update all blocks at each iteration. This variant is known as a *Jacobi* update. In the case of mixture of exponential family distribution, we have the following theorem.

**Theorem 14** *The Jacobi update is a natural-gradient update when the variational distribution is a member of the mixture of exponential family distribution.*

**Proof** When the variational distribution is a mixture of exponential family distribution, according to Lemma 13, the Jacobi update is

$$\boldsymbol{\lambda}_w^{k+1} = \boldsymbol{\lambda}_w^k + \beta \widehat{\nabla}_{m_w^k} \mathcal{L}(\boldsymbol{\lambda}) \quad (206)$$

$$\boldsymbol{\lambda}_z^{k+1} = \boldsymbol{\lambda}_z^k + \beta \widehat{\nabla}_{m_z^k} \mathcal{L}(\boldsymbol{\lambda}), \quad (207)$$

where  $\beta_w^k = \beta_z^k = \beta$ . It is easy to see that the Jacobi update is the natural-gradient update due to Theorem 1. ■

Note that Lemma 3 shows that the Fisher information matrix is block-diagonal in the case of mixture of exponential family distribution. We can generalize Theorem 15 and applied it to the multi-linear exponential family which gives the following theorem.

**Theorem 15** *The Jacobi update is a scaled gradient descent update where the scaling matrix is a block-diagonal approximated Fisher information matrix, if the variational distribution is either a mixture of exponential family distribution or a multi-linear exponential family distribution.*

**Proof** When the variational distribution is a mixture of exponential family distribution, by Theorem 14, the Jacobi update is a natural-gradient update. Recall that natural-gradient update is a scaled gradient update, where the scaling matrix is the Fisher information matrix. By Lemma 3, we know that the Fisher information matrix is block-diagonal.

Now, we consider the case when the variational distribution is a multi-linear exponential distribution. According to according to Lemma 13, we know the Jacobi update is

$$\boldsymbol{\lambda}_j^{k+1} = \boldsymbol{\lambda}_j^k + \beta \widehat{\nabla}_{m_j^k} \mathcal{L}(\boldsymbol{\lambda}) \quad j \in \{1, \dots, N\}. \quad (208)$$



Since  $\boldsymbol{\lambda}_{-j}^k$  is known, by the chain rule, we have

$$\widehat{\nabla}_{m_j^k} \widetilde{\mathcal{L}}(\boldsymbol{\lambda}) = \left( \widehat{\nabla}_{\boldsymbol{\lambda}_j^k} \mathbf{m}_j \right)^{-1} \widehat{\nabla}_{\boldsymbol{\lambda}_j^k} \mathcal{L}(\boldsymbol{\lambda}) \quad (209)$$

$$= \left( \widehat{\nabla}_{\boldsymbol{\lambda}_j^k} \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k)} \left[ \phi_j(\mathbf{z}, \boldsymbol{\lambda}_{-j}^k) \right] \right)^{-1} \widehat{\nabla}_{\boldsymbol{\lambda}_j^k} \mathcal{L}(\boldsymbol{\lambda}) \quad (210)$$

Recall that  $q(\mathbf{z}|\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k)$  is an one-parameter exponential family distribution given that  $\boldsymbol{\lambda}_{-j}$  is known. We have

$$\mathbf{0} = \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k)} \left[ \nabla_{\boldsymbol{\lambda}_j} \log q(\mathbf{z}|\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k) \right] \quad (211)$$

$$= \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k)} \left[ \nabla_{\boldsymbol{\lambda}_j} \left[ \langle \phi_j(\mathbf{z}, \boldsymbol{\lambda}_{-j}^k), \boldsymbol{\lambda}_j \rangle + r_j(\boldsymbol{\lambda}_{-j}^k) - A_z(\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k) \right] \right] \quad (212)$$

$$= \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k)} \left[ \phi_j(\mathbf{z}, \boldsymbol{\lambda}_{-j}^k) - \nabla_{\boldsymbol{\lambda}_j} A_z(\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k) \right] \quad (213)$$

$$= \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k)} \left[ \phi_j(\mathbf{z}, \boldsymbol{\lambda}_{-j}^k) \right] - \nabla_{\boldsymbol{\lambda}_j} A_z(\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k) \quad (214)$$

Let's define

$$\mathbf{F}_{\boldsymbol{\lambda}_j} = \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k)} \left[ \nabla_{\boldsymbol{\lambda}_j}^2 A_z(\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k) \right] = \nabla_{\boldsymbol{\lambda}_j}^2 A_z(\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k). \quad (215)$$

Due to Eq 214, we know that

$$\widehat{\nabla}_{\boldsymbol{\lambda}_j^k} \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda}_j, \boldsymbol{\lambda}_{-j}^k)} \left[ \phi_j(\mathbf{z}, \boldsymbol{\lambda}_{-j}^k) \right] = \mathbf{F}_{\boldsymbol{\lambda}_j^k}. \quad (216)$$

Therefore, the Jacobi update is

$$\boldsymbol{\lambda}_j^{k+1} = \boldsymbol{\lambda}_j^k + \beta \left( \mathbf{F}_{\boldsymbol{\lambda}_j^k} \right)^{-1} \widehat{\nabla}_{\boldsymbol{\lambda}_j^k} \mathcal{L}(\boldsymbol{\lambda}) \quad j \in \{1, \dots, N\}. \quad (217)$$

Recall that the Fisher information matrix of  $q(\mathbf{z}|\boldsymbol{\lambda})$  is

$$\mathbf{F}_{\boldsymbol{\lambda}^k} = - \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda}^k)} \left[ \widehat{\nabla}_{\boldsymbol{\lambda}^k}^2 \log q(\mathbf{z}|\boldsymbol{\lambda}) \right] \quad (218)$$

$$= - \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda}^k)} \begin{bmatrix} \widehat{\nabla}_{\boldsymbol{\lambda}_1^k}^2 \log q(\mathbf{z}|\boldsymbol{\lambda}) & \cdots & \widehat{\nabla}_{\boldsymbol{\lambda}_1^k} \widehat{\nabla}_{\boldsymbol{\lambda}_N^k} \log q(\mathbf{z}|\boldsymbol{\lambda}) \\ \cdots & \cdots & \cdots \\ \widehat{\nabla}_{\boldsymbol{\lambda}_N^k} \widehat{\nabla}_{\boldsymbol{\lambda}_1^k} \log q(\mathbf{z}|\boldsymbol{\lambda}) & \cdots & \widehat{\nabla}_{\boldsymbol{\lambda}_N^k}^2 \log q(\mathbf{z}|\boldsymbol{\lambda}) \end{bmatrix} \quad (219)$$

$$= - \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda}^k)} \begin{bmatrix} -\widehat{\nabla}_{\boldsymbol{\lambda}_1^k}^2 A_z(\boldsymbol{\lambda}) & \cdots & \widehat{\nabla}_{\boldsymbol{\lambda}_1^k} \widehat{\nabla}_{\boldsymbol{\lambda}_N^k} \log q(\mathbf{z}|\boldsymbol{\lambda}) \\ \cdots & \cdots & \cdots \\ \widehat{\nabla}_{\boldsymbol{\lambda}_N^k} \widehat{\nabla}_{\boldsymbol{\lambda}_1^k} \log q(\mathbf{z}|\boldsymbol{\lambda}) & \cdots & -\widehat{\nabla}_{\boldsymbol{\lambda}_N^k}^2 A_z(\boldsymbol{\lambda}) \end{bmatrix} \quad (220)$$

$$= \begin{bmatrix} \mathbf{F}_{\boldsymbol{\lambda}_1^k} & \cdots & -\widehat{\nabla}_{\boldsymbol{\lambda}_1^k} \widehat{\nabla}_{\boldsymbol{\lambda}_N^k} \log q(\mathbf{z}|\boldsymbol{\lambda}) \\ \cdots & \cdots & \cdots \\ -\widehat{\nabla}_{\boldsymbol{\lambda}_N^k} \widehat{\nabla}_{\boldsymbol{\lambda}_1^k} \log q(\mathbf{z}|\boldsymbol{\lambda}) & \cdots & \mathbf{F}_{\boldsymbol{\lambda}_N^k} \end{bmatrix} \quad (221)$$

From Eq 221, it is easy to see that the Jacobi update is a scaled gradient descent update where the scaling matrix at iteration  $k$  is the block approximation of the Fisher information

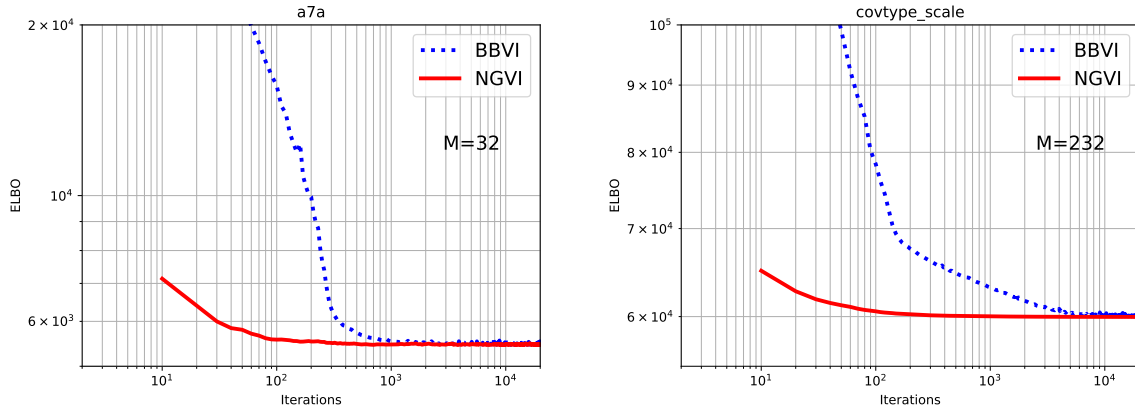


Figure 2: This figure demonstrates a fast convergence of NGVI over BBVI to approximate the posterior distribution of Bayesian logistic regression using a scale mixture of Gaussians with full covariance matrix. The plot shows the ELBO of the whole training set obtained using  $5 \times 10^4$  MC samples. For both algorithms, we used mini-batches by using 10 MC samples to compute stochastic approximations.  $M$  denotes that the size of a mini-batch. For BBVI, we use the Adam optimizer. We can see that our method is faster than BBVI using Adam.

matrix as shown below.

$$\widehat{\mathbf{F}}_{\lambda^k} := \begin{bmatrix} \mathbf{F}_{\lambda_1^k} & \cdots & \mathbf{0} \\ \cdots & \cdots & \cdots \\ \mathbf{0} & \cdots & \mathbf{F}_{\lambda_N^k} \end{bmatrix} \quad (222)$$

Finally, it is easy to check that the scaling matrix  $\widehat{\mathbf{F}}_{\lambda^k}$  is positive semi-definitive. If the minimality holds, it is positive definite.  $\blacksquare$

## Appendix G. Results for Multivariate t-Distribution

Figure 2 shows a fast convergence of NGVI over BBVI to approximate the posterior distribution of Bayesian logistic regression with a Student’s t prior expressed as a scale mixture of Gaussians discussed at Eq. 92. For this model, we use a scale mixture of Gaussians with full covariance matrix as the variational distribution.