

A. Proof of Theorem 2

In this section, we provide a proof for Theorem 2:

For an MCEF representation given in Definition 3, the FIM $\mathbf{F}_{wz}(\boldsymbol{\lambda})$ given in (9) is positive-definite and invertible for all $\boldsymbol{\lambda} \in \Omega$.	We prove this using a sequence of lemmas.
--	---

Lemma 1 $\log q(\mathbf{w}, \mathbf{z}|\boldsymbol{\lambda})$ is twice differentiable with respect to $\boldsymbol{\lambda}$.

Proof: From Definition 2, we see that the $\log q(\mathbf{z}, \mathbf{w})$ is differentiable when $A_w(\boldsymbol{\lambda}_w)$ and $A_z(\boldsymbol{\lambda}_z, \mathbf{w})$ are twice differentiable for each \mathbf{w} sampled from $q(\mathbf{w})$. Since $q(\mathbf{w})$ is an EF, $A_w(\boldsymbol{\lambda}_w)$ is twice differentiable (Johansen, 1979). Similarly, since conditioned on \mathbf{w} , $q(\mathbf{z}|\mathbf{w})$ is also an EF, $A_z(\boldsymbol{\lambda}_z, \mathbf{w})$ is twice differentiable too. Therefore, the log of the joint distribution is twice differentiable. \square

Lemma 2 The FIM $\mathbf{F}_{wz}(\boldsymbol{\lambda})$ is block-diagonal with two blocks:

$$\mathbf{F}_{wz}(\boldsymbol{\lambda}) = \begin{bmatrix} \mathbf{F}_z(\boldsymbol{\lambda}) & \mathbf{0} \\ \mathbf{0} & \mathbf{F}_w(\boldsymbol{\lambda}_w) \end{bmatrix}, \quad (32)$$

where $\mathbf{F}_w(\boldsymbol{\lambda}_w)$ is the FIM of $q(\mathbf{w})$ and $\mathbf{F}_z(\boldsymbol{\lambda})$ is the expected of the FIM of $q(\mathbf{z}|\mathbf{w})$ where expectation is taken under $q(\mathbf{w})$ as shown below:

$$\begin{aligned} \mathbf{F}_w(\boldsymbol{\lambda}_w) &:= -\mathbb{E}_{q(\mathbf{w})} [\nabla_{\boldsymbol{\lambda}_w}^2 \log q(\mathbf{w}|\boldsymbol{\lambda}_w)] \\ \mathbf{F}_z(\boldsymbol{\lambda}) &:= -\mathbb{E}_{q(\mathbf{w})} [\mathbb{E}_{q(\mathbf{z}|\mathbf{w})} [\nabla_{\boldsymbol{\lambda}_z}^2 \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)]] , \end{aligned}$$

Proof: By Lemma 1, $\log q(\mathbf{w}, \mathbf{z}|\boldsymbol{\lambda})$ is twice differentiable, so the FIM is well defined. Below, we simplify the FIM to show that it has a block-diagonal structure. The first step below follows from the definition of the FIM. The second step is simply writing the FIM in a 2×2 block corresponding to $\boldsymbol{\lambda}_z$ and $\boldsymbol{\lambda}_w$. In the third step, we write the joint as the product of $q(\mathbf{z}|\mathbf{w})$ and $q(\mathbf{w})$. The fourth step is obtained since the two blocks are separable in $\boldsymbol{\lambda}_z$ and $\boldsymbol{\lambda}_w$. In the fifth step, we take the expectation inside which give us the desired result in the last step.

$$\begin{aligned} \mathbf{F}_{wz}(\boldsymbol{\lambda}) &= -\mathbb{E}_{q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda})} [\nabla_{\boldsymbol{\lambda}}^2 \log q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda})] \\ &= -\mathbb{E}_{q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda})} \begin{bmatrix} \nabla_{\boldsymbol{\lambda}_z}^2 \log q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda}) & \nabla_{\boldsymbol{\lambda}_w} \nabla_{\boldsymbol{\lambda}_z^T} \log q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda}) \\ \nabla_{\boldsymbol{\lambda}_z} \nabla_{\boldsymbol{\lambda}_w^T} \log q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda}) & \nabla_{\boldsymbol{\lambda}_w}^2 \log q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda}) \end{bmatrix} \\ &= -\mathbb{E}_{q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda})} \begin{bmatrix} \nabla_{\boldsymbol{\lambda}_z}^2 (\log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) + \log q(\mathbf{w}|\boldsymbol{\lambda}_w)) & \nabla_{\boldsymbol{\lambda}_w} \nabla_{\boldsymbol{\lambda}_z^T} (\log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) + \log q(\mathbf{w}|\boldsymbol{\lambda}_w)) \\ \nabla_{\boldsymbol{\lambda}_z} \nabla_{\boldsymbol{\lambda}_w^T} (\log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) + \log q(\mathbf{w}|\boldsymbol{\lambda}_w)) & \nabla_{\boldsymbol{\lambda}_w}^2 (\log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) + \log q(\mathbf{w}|\boldsymbol{\lambda}_w)) \end{bmatrix} \\ &= -\mathbb{E}_{q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda})} \begin{bmatrix} \nabla_{\boldsymbol{\lambda}_z}^2 \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) & \mathbf{0} \\ \mathbf{0} & \nabla_{\boldsymbol{\lambda}_w}^2 \log q(\mathbf{w}|\boldsymbol{\lambda}_w) \end{bmatrix} \\ &= -\begin{bmatrix} \mathbb{E}_{q(\mathbf{z}, \mathbf{w}|\boldsymbol{\lambda})} [\nabla_{\boldsymbol{\lambda}_z}^2 \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)] & \mathbf{0} \\ \mathbf{0} & \mathbb{E}_{q(\mathbf{w}|\boldsymbol{\lambda}_w)} [\nabla_{\boldsymbol{\lambda}_w}^2 \log q(\mathbf{w}|\boldsymbol{\lambda}_w)] \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{F}_z(\boldsymbol{\lambda}) & \mathbf{0} \\ \mathbf{0} & \mathbf{F}_w(\boldsymbol{\lambda}_w) \end{bmatrix} \end{aligned}$$

\square

Lemma 3 The first block of the FIM matrix \mathbf{F}_z is equal to the derivative of the expectation parameter $\mathbf{m}_z(\boldsymbol{\lambda})$:

$$\mathbf{F}_z(\boldsymbol{\lambda}) := \nabla_{\boldsymbol{\lambda}_z} \mathbf{m}_z^T(\boldsymbol{\lambda})$$

Proof: We first show that the gradient of $A_z(\lambda_z, \mathbf{w})$ with respect to λ_z is equal to $\mathbb{E}_{q(z|\mathbf{w})} [\phi_z(\mathbf{z}, \mathbf{w})]$. By using the definition of $A_z(\lambda_z, \mathbf{w})$, this is straightforward to show:

$$\begin{aligned} \nabla_{\lambda_z} A_z(\lambda_z, \mathbf{w}) &= \nabla_{\lambda_z} \log \int h_z(\mathbf{z}, \mathbf{w}) \exp [\langle \phi_z(\mathbf{z}, \mathbf{w}), \lambda_z \rangle] d\mathbf{z} \\ &= \frac{\int \nabla_{\lambda_z} h_z(\mathbf{z}, \mathbf{w}) \exp [\langle \phi_z(\mathbf{z}, \mathbf{w}), \lambda_z \rangle] d\mathbf{z}}{\int h_z(\mathbf{z}, \mathbf{w}) \exp [\langle \phi_z(\mathbf{z}, \mathbf{w}), \lambda_z \rangle] d\mathbf{z}} \\ &= \frac{\int \phi_z(\mathbf{z}, \mathbf{w}) h_z(\mathbf{z}, \mathbf{w}) \exp [\langle \phi_z(\mathbf{z}, \mathbf{w}), \lambda_z \rangle] d\mathbf{z}}{\int h_z(\mathbf{z}, \mathbf{w}) \exp [\langle \phi_z(\mathbf{z}, \mathbf{w}), \lambda_z \rangle] d\mathbf{z}} \end{aligned} \quad (33)$$

$$= \mathbb{E}_{q(z|\mathbf{w})} [\phi_z(\mathbf{z}, \mathbf{w})] \quad (34)$$

Using this, the expectation parameter \mathbf{m}_z is simply the expected value of the gradient of the log-partition function.

$$\mathbf{m}_z = \mathbb{E}_{q(w)q(z|\mathbf{w})} [\phi_z(\mathbf{z}, \mathbf{w})] = \mathbb{E}_{q(w)} [\nabla_{\lambda_z} A_z(\lambda_z, \mathbf{w})] \quad (35)$$

Using this, it is easy to show the result by simply using the definition of the conditional EF, as shown below:

$$\begin{aligned} \mathbf{F}_z(\lambda) &= -\mathbb{E}_{q(z, \mathbf{w})} [\nabla_{\lambda_z}^2 \log q(\mathbf{z}|\mathbf{w}, \lambda_z)] \\ &= -\mathbb{E}_{q(z, \mathbf{w}|\lambda)} [\nabla_{\lambda_z}^2 (\log h_z(\mathbf{z}, \mathbf{w}) + \langle \phi_z(\mathbf{z}, \mathbf{w}), \lambda_z \rangle - A_z(\lambda_z, \mathbf{w}))] \\ &= \mathbb{E}_{q(z, \mathbf{w}|\lambda)} [\nabla_{\lambda_z}^2 A_z(\lambda_z, \mathbf{w})] \\ &= \mathbb{E}_{q(w|\lambda_w)} [\nabla_{\lambda_z}^2 A_z(\lambda_z, \mathbf{w})] \\ &= \nabla_{\lambda_z} \mathbb{E}_{q(w|\lambda_w)} [\nabla_{\lambda_z^T} A_z(\lambda_z, \mathbf{w})] \\ &= \nabla_{\lambda_z} \mathbf{m}_z^T \end{aligned}$$

□

Lemma 4 Let $\Omega_w \times \Omega_z$ be relatively open. If the mapping $\mathbf{m}_w(\cdot) : \Omega_w \rightarrow \mathcal{M}_w$ is one-to-one, and, given every $\lambda_w \in \Omega_w$, the conditional mapping $\mathbf{m}_z(\cdot, \lambda_w) : \Omega_z \rightarrow \mathcal{M}_z$ is one-to-one, then $\mathbf{F}_{wz}(\lambda)$ is positive-definite in $\Omega_z \times \Omega_w$.

Proof: When the mapping \mathbf{m}_w is one-to-one, $q(\mathbf{w}|\lambda_w)$ is a minimal EF, and given that Ω_w is relatively open, using the result discussed in Section 2, we conclude that the second block $\mathbf{F}_w(\lambda_w)$ of $\mathbf{F}_{wz}(\lambda)$ given in (32) is positive definite and invertible for all Ω_z . Now we prove that the first block $\mathbf{F}_z(\lambda)$ is also positive definite.

The steps below establish the positive-semi definiteness first. The first step is simply the definition of the FIM, while the second step is obtained by using the fact that $\nabla \log f(\lambda) = \nabla f(\lambda)/f(\lambda)$. The third step is obtained by using the chain-rule, and the fourth step simply uses the log-trick above to simplify the second term. In the fifth step, we take the derivative out of the first term which cancels out $q(\mathbf{z}|\mathbf{w}, \lambda_z)$. The last step is straightforward since the outer products are always nonnegative.

$$\begin{aligned} \nabla_{\lambda_z}^2 A_z(\lambda_z, \mathbf{w}) &= -\mathbb{E}_{q(z|\mathbf{w})} [\nabla_{\lambda_z}^2 \log q(\mathbf{z}|\mathbf{w}, \lambda_z)], \\ &= -\mathbb{E}_{q(z|\mathbf{w})} \left[\nabla_{\lambda_z} \left(\frac{\nabla_{\lambda_z^T} q(\mathbf{z}|\mathbf{w}, \lambda_z)}{q(\mathbf{z}|\mathbf{w}, \lambda_z)} \right) \right], \\ &= -\mathbb{E}_{q(z|\mathbf{w})} \left[\frac{\nabla_{\lambda_z}^2 q(\mathbf{z}|\mathbf{w}, \lambda_z)}{q(\mathbf{z}|\mathbf{w}, \lambda_z)} - \frac{\nabla_{\lambda_z} q(\mathbf{z}|\mathbf{w}, \lambda_z)}{q(\mathbf{z}|\mathbf{w}, \lambda_z)} \frac{\nabla_{\lambda_z^T} q(\mathbf{z}|\mathbf{w}, \lambda_z)}{q(\mathbf{z}|\mathbf{w}, \lambda_z)} \right] \\ &= \mathbb{E}_{q(z|\mathbf{w})} \left[-\frac{\nabla_{\lambda_z}^2 q(\mathbf{z}|\mathbf{w}, \lambda_z)}{q(\mathbf{z}|\mathbf{w}, \lambda_z)} \right] + \mathbb{E}_{q(z|\mathbf{w})} [\nabla_{\lambda_z} \log q(\mathbf{z}|\mathbf{w}, \lambda_z) \nabla_{\lambda_z^T} \log q(\mathbf{z}|\mathbf{w}, \lambda_z)], \\ &= \int -\nabla_{\lambda_z}^2 q(\mathbf{z}|\mathbf{w}, \lambda_z) d\mathbf{z} + \mathbb{E}_{q(z|\mathbf{w})} [\nabla_{\lambda_z} \log q(\mathbf{z}|\mathbf{w}, \lambda_z) \nabla_{\lambda_z^T} \log q(\mathbf{z}|\mathbf{w}, \lambda_z)], \\ &= \underbrace{-\nabla_{\lambda_z}^2 \int q(\mathbf{z}|\mathbf{w}, \lambda_z) d\mathbf{z}}_{=0} + \mathbb{E}_{q(z|\mathbf{w})} [\nabla_{\lambda_z} \log q(\mathbf{z}|\mathbf{w}, \lambda_z) \nabla_{\lambda_z^T} \log q(\mathbf{z}|\mathbf{w}, \lambda_z)], \\ &= \mathbb{E}_{q(z|\mathbf{w})} [\nabla_{\lambda_z} \log q(\mathbf{z}|\mathbf{w}, \lambda_z) \nabla_{\lambda_z^T} \log q(\mathbf{z}|\mathbf{w}, \lambda_z)] \succeq \mathbf{0}. \end{aligned} \quad (36)$$

Using Lemma 3 and (35), we see that FIM is positive semi-definite:

$$\mathbf{F}_z(\boldsymbol{\lambda}) = \nabla_{\lambda_z} \mathbf{m}_z^T = \nabla_{\lambda_z} \mathbb{E}_{q(w|\lambda_w)} [\nabla_{\lambda_z^T} A_z(\boldsymbol{\lambda}_z, \mathbf{w})] = \mathbb{E}_{q(w)} [\nabla_{\lambda_z^T}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w})] \succeq \mathbf{0}$$

Now, we prove the final claim that, for every $\boldsymbol{\lambda}_w \in \Omega_w$, if the conditional mapping $\mathbf{m}_z(\cdot, \boldsymbol{\lambda}_w)$ is one-to-one, then $\mathbf{F}_z(\boldsymbol{\lambda})$ is positive definite. We will prove this statement by contradiction. Suppose there exists $\boldsymbol{\lambda}$ such that $\mathbf{F}_z(\boldsymbol{\lambda})$ is positive semi-definite, since $\mathbf{F}_z(\boldsymbol{\lambda})$ is positive semi-definite, there exists a non-zero vector \mathbf{a} such that $\mathbf{a}^T \mathbf{F}_z(\boldsymbol{\lambda}) \mathbf{a} = 0$. Simplifying below, we show that this leads to a contradiction. The first and second step are obtained by simply plugging (36), while the third step is obtained by using the definition of $q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z)$ and the fourth step is obtained by using (34). The last step is obtained by noting that the quantity is simply the variance of $\mathbf{a}^T \phi_z(\mathbf{z}, \mathbf{w})$ conditioned on \mathbf{w} .

$$\begin{aligned} \mathbf{a}^T \mathbf{F}_z(\boldsymbol{\lambda}) \mathbf{a} &= \mathbb{E}_{q(w)} [\mathbf{a}^T \nabla_{\lambda_z}^2 A_z(\boldsymbol{\lambda}_z, \mathbf{w}) \mathbf{a}] \\ &= \mathbb{E}_{q(w)} [\mathbf{a}^T \mathbb{E}_{q(z|\mathbf{w})} \{ \nabla_{\lambda_z} \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) \nabla_{\lambda_z^T} \log q(\mathbf{z}|\mathbf{w}, \boldsymbol{\lambda}_z) \} \mathbf{a}] \\ &= \mathbb{E}_{q(w)q(z|\mathbf{w})} [\mathbf{a}^T (\phi_z(\mathbf{z}, \mathbf{w}) - \nabla_{\lambda_z} A_z(\boldsymbol{\lambda}_z, \mathbf{w})) (\phi_z(\mathbf{z}, \mathbf{w}) - \nabla_{\lambda_z} A_z(\boldsymbol{\lambda}_z, \mathbf{w}))^T \mathbf{a}] \\ &= \mathbb{E}_{q(w)q(z|\mathbf{w})} [\mathbf{a}^T (\phi_z(\mathbf{z}, \mathbf{w}) - \mathbb{E}_{q(z|\mathbf{w})} [\phi_z(\mathbf{z}, \mathbf{w})]) (\phi_z(\mathbf{z}, \mathbf{w}) - \mathbb{E}_{q(z|\mathbf{w})} [\phi_z(\mathbf{z}, \mathbf{w})])^T \mathbf{a}] \\ &= \mathbb{E}_{q(w)} \mathbb{V}_{q(z|\mathbf{w})} [\mathbf{a}^T \phi_z(\mathbf{z}, \mathbf{w})] \end{aligned}$$

The expectation of a function positive quantity is equal to zero only when each function value is equal to zero, therefore for the above to be zeros, we need $\mathbf{a}^T \phi_z(\mathbf{z}, \mathbf{w}) = 0$. However, as we show below, this is not possible since the representation $q(\mathbf{z}|\mathbf{w})$ is minimal conditioned on \mathbf{w} .

Since Ω_z is relatively open, there exists a small $\delta > 0$ to always be able to obtain a perturbed version $\boldsymbol{\lambda}'_z = \boldsymbol{\lambda}_z + \delta \mathbf{a}$, such that $\boldsymbol{\lambda}'_z \in \Omega_z$. Since the conditional mapping is one-to-one, $\mathbf{m}_z(\boldsymbol{\lambda}'_z, \boldsymbol{\lambda}_w) \neq \mathbf{m}_z(\boldsymbol{\lambda}_z, \boldsymbol{\lambda}_w)$. By using (35) and (33), when $\mathbf{a}^T \phi_z(\mathbf{z}, \mathbf{w}) = 0$, we get a contradiction:

$$\begin{aligned} \mathbf{m}_z(\boldsymbol{\lambda}'_z, \boldsymbol{\lambda}_w) &= \mathbb{E}_{q(w|\lambda_w)} [\nabla_{\lambda'_z} A_z(\boldsymbol{\lambda}'_z, \mathbf{w})] \\ &= \mathbb{E}_{q(w|\lambda_w)} \left[\frac{\int \phi_z(\mathbf{z}, \mathbf{w}) h_z(\mathbf{z}, \mathbf{w}) \exp [\langle \phi_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}'_z \rangle] d\mathbf{z}}{\int h_z(\mathbf{z}, \mathbf{w}) \exp [\langle \phi_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}'_z \rangle] d\mathbf{z}} \right] \end{aligned} \quad (37)$$

$$\begin{aligned} &= \mathbb{E}_{q(w|\lambda_w)} \left[\frac{\int \phi_z(\mathbf{z}, \mathbf{w}) h_z(\mathbf{z}, \mathbf{w}) \exp [\langle \phi_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}_z \rangle] d\mathbf{z}}{\int h_z(\mathbf{z}, \mathbf{w}) \exp [\langle \phi_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}_z \rangle] d\mathbf{z}} \right] \\ &= \mathbf{m}_z(\boldsymbol{\lambda}_z, \boldsymbol{\lambda}_w) \end{aligned} \quad (38)$$

where we can move from (37) to (38), since $\langle \phi_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}'_z \rangle = \langle \phi_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}_z + \delta \mathbf{a} \rangle = \langle \phi_z(\mathbf{z}, \mathbf{w}), \boldsymbol{\lambda}_z \rangle$. Due to the contradiction, $\mathbf{F}_z(\boldsymbol{\lambda})$ must be positive definite. This proves that both the blocks are positive definite and invertible. \square

Lemma 5 *The gradient with respect to $\boldsymbol{\lambda}$ can be expressed as the gradient with respect to \mathbf{m} :*

$$\nabla_{\boldsymbol{\lambda}} \mathcal{L} = [\nabla_{\boldsymbol{\lambda}} \mathbf{m}^T] \nabla_{\mathbf{m}} \mathcal{L} = [\mathbf{F}_{wz}(\boldsymbol{\lambda})] \nabla_{\mathbf{m}} \mathcal{L} \quad (39)$$

Proof: Using Lemma 3, and chain rule, we can establish the results for $\boldsymbol{\lambda}_z$:

$$\nabla_{\lambda_z} \mathcal{L} = \nabla_{\lambda_z} \mathbf{m}_z^T(\boldsymbol{\lambda}) [\nabla_{\mathbf{m}_z} \mathcal{L}] = \mathbf{F}_z(\boldsymbol{\lambda}) [\nabla_{\mathbf{m}_z} \mathcal{L}]$$

For $\boldsymbol{\lambda}_w$, this result holds trivially, which proves the statement. \square

B. Finite Mixture of Gaussians

The finite mixture of EF distribution has the following conditional distribution $q(\mathbf{z}|w)$:

$$\begin{aligned} q(\mathbf{z}|w) &= \sum_{c=1}^K \mathbb{I}_c(w) q(\mathbf{z}|\boldsymbol{\lambda}_c) = \sum_{c=1}^K \mathbb{I}_c(w) h_z(\mathbf{z}) \exp [\langle \boldsymbol{\lambda}_c, \boldsymbol{\phi}_z(\mathbf{z}) \rangle - A_z(\boldsymbol{\lambda}_c)] \\ &= h_z(\mathbf{z}) \exp \left\{ \sum_{c=1}^K \langle \mathbb{I}_c(w) \boldsymbol{\phi}_z(\mathbf{z}), \boldsymbol{\lambda}_c \rangle - \sum_{c'=1}^K \mathbb{I}_{c'}(w) A_z(\boldsymbol{\lambda}_{c'}) \right\} \end{aligned}$$

where we assume each component admits the same parametric form.

From the above expression and using the EF form for the multinomial distribution, we can write the sufficient statistics, natural parameters, and expectation parameters as shown below, where $\mathbf{m}_c := \mathbb{E}_{q(\mathbf{z}|w=c)} [\boldsymbol{\phi}_z(\mathbf{z})]$ is the expectation parameter of a component $q(\mathbf{z}|w=c)$.

$$\begin{bmatrix} \mathbb{I}_1(w) \\ \mathbb{I}_2(w) \\ \vdots \\ \mathbb{I}_{K-1}(w) \\ \mathbb{I}_1(w) \boldsymbol{\phi}_z(\mathbf{z}) \\ \mathbb{I}_2(w) \boldsymbol{\phi}_z(\mathbf{z}) \\ \vdots \\ \mathbb{I}_K(w) \boldsymbol{\phi}_z(\mathbf{z}) \end{bmatrix} \quad \begin{bmatrix} \log(\pi_1/\pi_K) \\ \log(\pi_2/\pi_K) \\ \vdots \\ \log(\pi_{K-1}/\pi_K) \\ \boldsymbol{\lambda}_1 \\ \boldsymbol{\lambda}_2 \\ \vdots \\ \boldsymbol{\lambda}_K \end{bmatrix} \quad \begin{bmatrix} \pi_1 \\ \pi_2 \\ \vdots \\ \pi_{K-1} \\ \pi_1 \mathbf{m}_1 \\ \pi_2 \mathbf{m}_2 \\ \vdots \\ \pi_K \mathbf{m}_K \end{bmatrix}$$

From the last two vectors, we can see that the mapping between $\boldsymbol{\lambda}$ and \mathbf{m} is going to be one-to-one, when each EF $q(\mathbf{z}|\boldsymbol{\lambda}_c)$ is minimal (which makes sure that mapping $\boldsymbol{\lambda}_c$ and \mathbf{m}_c is one-to-one), and all $\boldsymbol{\lambda}_c$ are distinct.

B.1. The Model and ELBO

We consider the following model: $p(\mathcal{D}, \mathbf{z}) = \prod_{n=1}^N p(\mathcal{D}_n|\mathbf{z})p(\mathbf{z})$. We approximate the posterior by using the finite mixture of EFs whose marginal is denoted by $q(\mathbf{z})$ as given in (15). The variational lower bound is given by the following:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\lambda}) &= \left[\sum_{n=1}^N [\log p(\mathcal{D}_n|\mathbf{z})] + \log \frac{p(\mathbf{z})}{q(\mathbf{z})} \right] \\ &= \mathbb{E}_{q(\mathbf{z})} [-h(\mathbf{z})], \text{ where } h(\mathbf{z}) := - \left[\log \frac{p(\mathbf{z})}{q(\mathbf{z})} + \sum_n \log p(\mathcal{D}_n|\mathbf{z}) \right]. \end{aligned}$$

Note that the lower bound is defined with the marginal $q(\mathbf{z})$ and the variable w is not part of the model but only the variational approximation $q(\mathbf{z}, w)$.

B.2. Finite Mixture of Gaussians Approximation

We now give details about the NGD update for finite mixture of Gaussians. Note that the NGD update for $\boldsymbol{\lambda}_z$ and $\boldsymbol{\lambda}_w$ can be computed separately since the FIM is block-diagonal. We first derive the NGD update for each component $q(\mathbf{z}|w=c)$, and then give an update for $\boldsymbol{\lambda}_w$.

As shown in Table 1, the natural and expectation parameters of the c 'th component is given as follows:

$$\begin{aligned} \boldsymbol{\Lambda}_c &:= -\frac{1}{2} \boldsymbol{\Sigma}_c^{-1} & \mathbf{M}_c &:= \pi_c (\boldsymbol{\mu}_c \boldsymbol{\mu}_c^T + \boldsymbol{\Sigma}_c) \\ \boldsymbol{\lambda}_c &:= \boldsymbol{\Sigma}_c^{-1} \boldsymbol{\mu}_c & \mathbf{m}_c &:= \pi_c \boldsymbol{\mu}_c \end{aligned}$$

The expectation parameters \mathbf{m}_c and \mathbf{M}_c are functions of $\pi_c, \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c$ and its gradient can be obtained in terms of the gradient

with respect to these quantities by using the chain rule. The final expressions are shown below:

$$\begin{aligned}\nabla_{\mu_c} \mathcal{L} &= \frac{1}{\pi_c} (\nabla_{\mu_c} \mathcal{L} - 2 [\nabla_{\Sigma_c} \mathcal{L}] \mu_c) \\ \nabla_{\Sigma_c} \mathcal{L} &= \frac{1}{\pi_c} (\nabla_{\Sigma_c} \mathcal{L})\end{aligned}$$

We can compute gradients with respect to μ_c and Σ_c by using the gradient and Hessian of $h(\mathbf{z})$ at a sample \mathbf{z} from $q(\mathbf{z}, w)$. This can be done by using the Bonnet's and Price's theorems (Bonnet, 1964; Price, 1958; Opper & Archambeau, 2009; Rezende et al., 2014). (Staines & Barber, 2012) and Lin et al. (2019) discuss the conditions of the target function $h(\mathbf{z})$ when it comes to applying these theorems. Firstly, we define $\delta_c = q(\mathbf{z}|w=c)/q(\mathbf{z}) := \mathcal{N}(\mathbf{z}|\mu_c, \Sigma_c) / \sum_{c'=1}^K \pi_{c'} \mathcal{N}(\mathbf{z}|\mu_{c'}, \Sigma_{c'})$. Using these theorems, we obtain the following stochastic-gradient estimations for the mixture of Gaussian case:

$$\begin{aligned}\nabla_{\mu_c} \mathcal{L}(\lambda) &= -\mathbb{E}_{q(\mathbf{z})} [q(w=c|\mathbf{z}) \nabla_{\mathbf{z}} h(\mathbf{z})] \approx -\pi_c \delta_c \nabla_{\mathbf{z}} h(\mathbf{z}) \\ \nabla_{\Sigma_c} \mathcal{L}(\lambda) &= -\mathbb{E}_{q(\mathbf{z})} [q(w=c|\mathbf{z}) \nabla_{\mathbf{z}}^2 h(\mathbf{z})] \approx -\frac{\pi_c \delta_c}{2} \nabla_{\mathbf{z}}^2 h(\mathbf{z}).\end{aligned}$$

where $q(w=c|\mathbf{z}) = \pi_c \delta_c$ and \mathbf{z} is sampled from $q(\mathbf{z})$.

We can then plug these gradient estimations in the natural-gradient update for Λ_c :

$$-\frac{1}{2} [\Sigma_c^{(\text{new})}]^{-1} \leftarrow -\frac{1}{2} \Sigma_c^{-1} + \beta \nabla_{\Sigma_c} \mathcal{L} \quad \Rightarrow \quad [\Sigma_c^{(\text{new})}]^{-1} \leftarrow \Sigma_c^{-1} + \beta \delta_c \nabla_{\mathbf{z}}^2 h(\mathbf{z})$$

Similar for λ_c :

$$\begin{aligned}[\Sigma_c^{(\text{new})}]^{-1} \mu_c^{(\text{new})} &\leftarrow \Sigma_c^{-1} \mu_c + \beta \nabla_{\mu_c} \mathcal{L} \\ &\leftarrow \Sigma_c^{-1} \mu_c + \beta \frac{1}{\pi_c} (\nabla_{\mu_c} \mathcal{L} - 2 [\nabla_{\Sigma_c} \mathcal{L}] \mu_c) \\ &\leftarrow \left[\Sigma_c^{-1} - 2 \frac{\beta}{\pi_c} [\nabla_{\Sigma_c} \mathcal{L}] \right] \mu_c + \beta \frac{1}{\pi_c} (\nabla_{\mu_c} \mathcal{L}) \\ &\leftarrow [\Sigma_c^{-1} + \beta \delta_c [\nabla_{\mathbf{z}}^2 h(\mathbf{z})]] \mu_c + \beta \frac{1}{\pi_c} (\nabla_{\mu_c} \mathcal{L}) \\ &\leftarrow [\Sigma_c^{(\text{new})}]^{-1} \mu_c + \beta \frac{1}{\pi_c} (\nabla_{\mu_c} \mathcal{L})\end{aligned}$$

This gives the following update (by using the stochastic gradients):

$$\mu_c^{(\text{new})} \leftarrow \mu_c - \beta \delta_c \Sigma_c^{(\text{new})} \nabla_{\mathbf{z}} h(\mathbf{z})$$

B.3. Natural Gradients for $q(w)$

Now, we give the update for $q(w|\lambda_w)$. Its natural parameter and expectation parameter are

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

To derive the gradients, we note that only $q(\mathbf{z})$ depends on π since the model does not contain this as a parameter. Therefore, we need the gradient of the variational approximation which can be written as follows:

$$\nabla_{\pi_c} q(\mathbf{z}) = \nabla_{\pi_c} \sum_{k=1}^K \pi_k q(\mathbf{z}|\lambda_k) = q(\mathbf{z}|\lambda_c) - q(\mathbf{z}|\lambda_K).$$

The second term appears because the last π_K depends on π_c .

For the convenience of our derivation, we will separate the lower bound into terms that depend on $q(\mathbf{z})$ and the rest of the terms, as shown below:

$$\begin{aligned}
 \nabla_{\pi_c} \mathcal{L}(\boldsymbol{\lambda}) &= \nabla_{\pi_c} \mathbb{E}_{q(\mathbf{z})} \left[\sum_{n=1}^N \log p(\mathcal{D}_n | \mathbf{z}) + \log p(\mathbf{z}) - \log q(\mathbf{z}) \right] \\
 &= \int \underbrace{\nabla_{\pi_c} q(\mathbf{z})}_{q(\mathbf{z} | \boldsymbol{\lambda}_c) - q(\mathbf{z} | \boldsymbol{\lambda}_K)} \left[\sum_{n=1}^N \log p(\mathcal{D}_n | \mathbf{z}) + \log p(\mathbf{z}) - \log q(\mathbf{z}) \right] d\mathbf{z} - \underbrace{\int q(\mathbf{z}) \nabla_{\pi_c} \log q(\mathbf{z}) d\mathbf{z}}_0 \\
 &= \int [q(\mathbf{z} | \boldsymbol{\lambda}_c) - q(\mathbf{z} | \boldsymbol{\lambda}_K)] \left[\sum_{n=1}^N \left(\log p(\mathcal{D}_n | \mathbf{z}) + \log \frac{p(\mathbf{z})}{q(\mathbf{z})} \right) \right] d\mathbf{z} \\
 &= \int q(\mathbf{z}) \left[\frac{q(\mathbf{z} | \boldsymbol{\lambda}_c)}{q(\mathbf{z})} - \frac{q(\mathbf{z} | \boldsymbol{\lambda}_K)}{q(\mathbf{z})} \right] \left[\sum_{n=1}^N \left(\log p(\mathcal{D}_n | \mathbf{z}) + \log \frac{p(\mathbf{z})}{q(\mathbf{z})} \right) \right] d\mathbf{z} \\
 &= \mathbb{E}_{q(\mathbf{z})} \left[(\delta_c - \delta_K) \underbrace{\left[\sum_{n=1}^N \left(\log p(\mathcal{D}_n | \mathbf{z}) + \log \frac{p(\mathbf{z})}{q(\mathbf{z})} \right) \right]}_{-h(\mathbf{z})} \right] \\
 &\approx -(\delta_c - \delta_K) h(\mathbf{z})
 \end{aligned}$$

where \mathbf{z} is a sample from $q(\mathbf{z})$.

Using this, we can perform the following NGD update:

$$\log(\pi_c / \pi_K) \leftarrow \log(\pi_c / \pi_K) - \beta(\delta_c - \delta_K) h(\mathbf{z})$$

B.4. Extension to Finite Mixture of EFs

The algorithm presented in Section 4.1 can be extended to handle generic minimal EF components. We now present a general gradient estimator to do so. The update of π_c remains unaltered, so we do not discuss them here. We only discuss how to update natural parameters $\boldsymbol{\lambda}_c$ of $q(\mathbf{z} | \boldsymbol{\lambda}_c)$.

The natural parameter and sufficient statistics are $\boldsymbol{\lambda}_c$ and $\mathbb{I}_c(w) \phi_z(\mathbf{z})$ respectively. We wish to perform the following update: $\boldsymbol{\lambda}_c \leftarrow \boldsymbol{\lambda}_c + \beta_z \nabla_{m_c} \mathcal{L}(\boldsymbol{\lambda})$. In general, we can compute the gradient $\nabla_{m_c} \mathcal{L}(\boldsymbol{\lambda})$ by computing the FIM of each component as shown below:

$$\nabla_{m_c} \mathcal{L}(\boldsymbol{\lambda}) = (\nabla_{\boldsymbol{\lambda}_c} \mathbf{m}_c)^{-1} \nabla_{\boldsymbol{\lambda}_c} \mathcal{L}(\boldsymbol{\lambda}) = (\nabla_{\boldsymbol{\lambda}_c} \mathbb{E}_{q(w, \mathbf{z})} [\mathbb{I}_c(w) \phi_z(\mathbf{z})])^{-1} \nabla_{\boldsymbol{\lambda}_c} \mathcal{L}(\boldsymbol{\lambda}),$$

Both of these gradients can be obtained given $\nabla_{\boldsymbol{\lambda}_c} \mathbf{z}$, where \mathbf{z} is a sample from $q(\mathbf{z})$ as shown below.

$$\begin{aligned}
 \nabla_{\boldsymbol{\lambda}_c} \mathbb{E}_{q(w, \mathbf{z})} [\mathbb{I}_c(w) \phi_z(\mathbf{z})] &= \nabla_{\boldsymbol{\lambda}_c} \int \pi_c q(\mathbf{z} | w = c) \phi_z(\mathbf{z}) d\mathbf{z} = \int \pi_c q(\mathbf{z} | w = c) \nabla_z [\phi_z(\mathbf{z})] [\nabla_{\boldsymbol{\lambda}_c} \mathbf{z}] d\mathbf{z} \\
 \nabla_{\boldsymbol{\lambda}_c} \mathcal{L}(\boldsymbol{\lambda}) &= \int \pi_c \nabla_{\boldsymbol{\lambda}_c} q(\mathbf{z} | w = c) [-h(\mathbf{z})] d\mathbf{z} + \underbrace{\int q(\mathbf{z}) [\nabla_{\boldsymbol{\lambda}_c} \log q(\mathbf{z})] d\mathbf{z}}_0 \\
 &= \int \pi_c q(\mathbf{z} | w = c) [\nabla_z (-h(\mathbf{z}))] [\nabla_{\boldsymbol{\lambda}_c} \mathbf{z}] d\mathbf{z}
 \end{aligned}$$

If we assume that $q(z|w)$ is an univariate continuous exponential family distribution, we can use the implicitly re-parameterization trick (Salimans & Knowles, 2013; Figurnov et al., 2018) to get:

$$\nabla_{\boldsymbol{\lambda}_c} \mathbf{z} = -\frac{\nabla_{\boldsymbol{\lambda}_c} Q_c(\mathbf{z} | \boldsymbol{\lambda}_c)}{q(\mathbf{z} | w = c)},$$

where $Q_c(\cdot|\lambda_c)$ is the cumulative distribution function (CDF) of $q(z|w=c)$. Therefore, we now can compute the required gradient as below:

$$\begin{aligned}\nabla_{\lambda_c} \mathcal{L}(\lambda) &= \int \pi_c q(z|w=c) [-\nabla_z h(z)] \left[-\frac{\nabla_{\lambda_c} Q_c(z|\lambda_c)}{q(z|w=c)} \right] dz \\ &= \mathbb{E}_{q(z)} \left[\pi_c \frac{q(z|w=c)}{q(z)} [-\nabla_z h(z)] \left[-\frac{\nabla_{\lambda_c} Q_c(z|\lambda_c)}{q(z|w=c)} \right] \right] \\ &= \mathbb{E}_{q(z)} \left[\frac{\pi_c \nabla_{\lambda_c} Q_c(z|\lambda_c)}{q(z)} [\nabla_z h(z)] \right],\end{aligned}$$

and

$$\begin{aligned}\nabla_{\lambda_c} \mathbb{E}_{q(w,z)} [\mathbb{I}_c(w) \phi_z(z)] &= \int \pi_c q(z|w=c) \nabla_z [\phi_z(z)] \left[-\frac{\nabla_{\lambda_c} Q_c(z|\lambda_c)}{q(z|w=c)} \right] dz \\ &= \mathbb{E}_{q(z)} \left[\frac{-\pi_c \nabla_{\lambda_c} Q_c(z|\lambda_c)}{q(z)} \nabla_z [\phi_z(z)] \right].\end{aligned}$$

This is not the most efficient way to compute NGDs, however, for the specific cases (e.g., Gaussian, exponential distribution, inverse Gaussian), we can get simplifications whenever gradient with respect to the expectation parameters are easy to compute.

The Birnbaum-Saunders distribution, which is a finite mixture of inverse Gaussians, is presented in Appendix C. This example is different from examples given in this section since we allow each mixing component takes a different and tied parametric form.

B.5. Result for the Toy Example

See Figure 6

C. Birnbaum-Saunders Distribution

Firstly, we denote the inverse Gaussian distribution by $\text{InvGauss}(z; \mu, v) = \left(\frac{v}{2\pi z^3}\right)^{1/2} \exp\left[-\frac{vz}{2\mu^2} - \frac{v}{2z} + \frac{v}{\mu}\right]$, where $z > 0$, $v > 0$, and $\mu > 0$. We consider the following mixture distribution.

$$\begin{aligned}q(w) &= p^{\mathbb{I}(w=0)}(1-p)^{\mathbb{I}(w=1)} \\ q(z|w) &= \mathbb{I}(w=0) \text{InvGauss}(z; \mu, v) + \mathbb{I}(w=1) \frac{z \text{InvGauss}(z; \mu, v)}{\mu},\end{aligned}\tag{40}$$

where $\frac{z}{\mu} \text{InvGauss}(z; \mu, v)$ is a normalized distribution since it is the distribution of $z = y^{-1}$ where y is distributed by $\text{InvGauss}(y; \mu^{-1}, v/\mu^2)$.

As we can observe from Eq (40), each mixing component has a distinct parametric form and variational parameters are shared between the components, which is different from examples discussed in Appendix B. According to Desmond (1986); Jørgensen et al. (1991); Balakrishnan & Kundu (2019), the marginal distribution is known as the Birnbaum-Saunders distribution (Birnbaum & Saunders, 1969) shown as below, where $p = \frac{1}{2}$.

$$\begin{aligned}q(z|v, \mu) &= \sum_w q(w) q(z|w) \\ &= \frac{1}{2} \left\{ \left(\frac{v}{2\pi z^3}\right)^{1/2} \exp\left[-\frac{vz}{2\mu^2} - \frac{v}{2z} + \frac{v}{\mu}\right] \left(1 + \frac{z}{\mu}\right) \right\} \\ &= \frac{\sqrt{v}}{2\sqrt{2\pi}} \left[\frac{1}{z^{3/2}} + \frac{1}{\mu z^{1/2}} \right] \exp\left[-\frac{vz}{2\mu^2} - \frac{v}{2z} + \frac{v}{\mu}\right] \\ &= \frac{\sqrt{v}}{2\mu\sqrt{2\pi}\mu} \left[\left(\frac{\mu}{z}\right)^{1/2} + \left(\frac{\mu}{z}\right)^{3/2} \right] \exp\left\{-\frac{v\left(\frac{z}{\mu} + \frac{\mu}{z} - 2\right)}{2\mu}\right\}\end{aligned}$$

Lemma 6 *The joint distribution of the Birnbaum-Saunders distribution given in (40) can be written in a conditional EF form.*

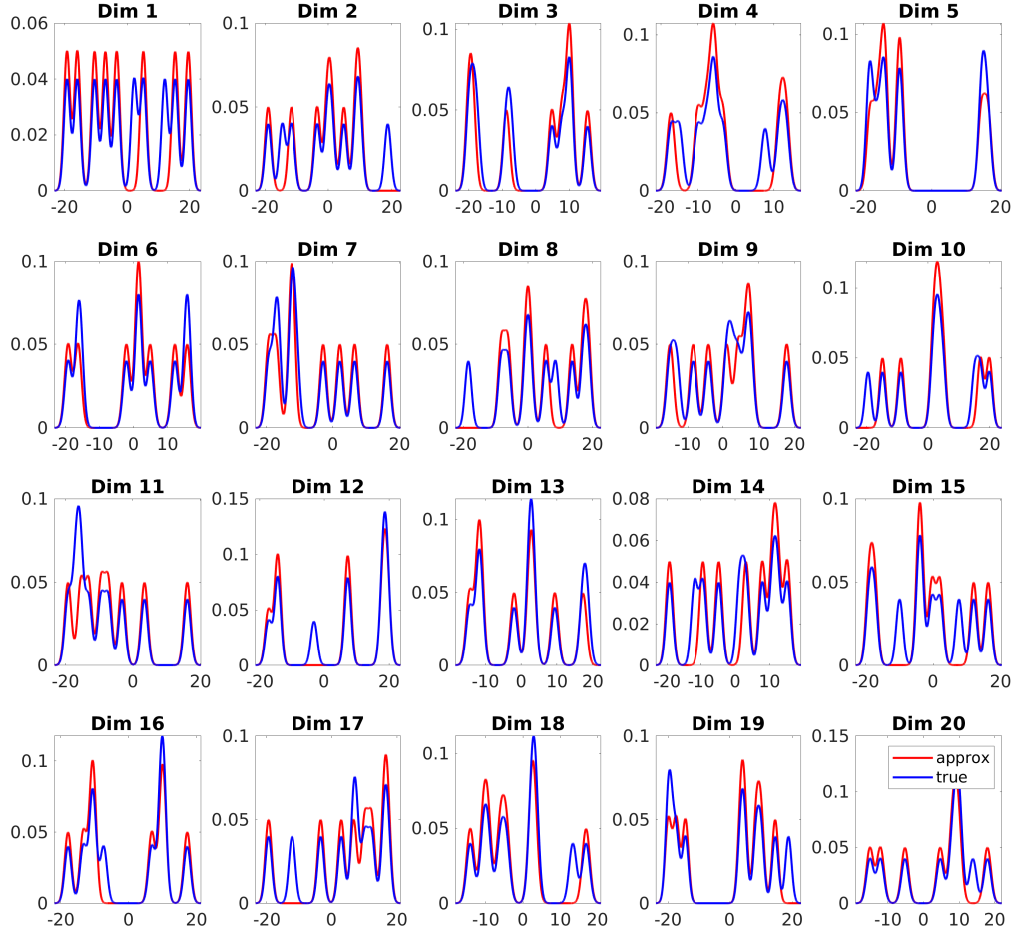


Figure 6. This is a complete version of the leftmost figure in Figure 2. The figure shows MOG approximation (with $K = 20$) to fit an MOG model with 10 components in a 20 dimensional problem.

Proof: It is obvious that $q(w)$ is Bernoulli distribution with $p = \frac{1}{2}$, which is an EF distribution. Now, we show that $q(z|w)$ is a conditional EF distribution as below.

$$\begin{aligned}
 q(z|w) &= \exp \left\{ \mathbb{I}(w=0) \left[-\frac{vz}{2\mu^2} - \frac{v}{2z} + \frac{v}{\mu} + \frac{1}{2} \log \frac{v}{2\pi z^3} \right] + \mathbb{I}(w=1) \left[-\frac{vz}{2\mu^2} - \frac{v}{2z} + \frac{v}{\mu} + \frac{1}{2} \log \frac{v}{2\pi z} - \log \mu \right] \right\} \\
 &= \exp \left\{ -\frac{vz}{2\mu^2} - \frac{v}{2z} + \frac{v}{\mu} + \mathbb{I}(w=0) \left[\frac{1}{2} \log \frac{v}{2\pi z^3} \right] + \mathbb{I}(w=1) \left[\frac{1}{2} \log \frac{v}{2\pi z} - \log \mu \right] \right\} \\
 &= \frac{1}{\sqrt{2\pi}} z^{-3\mathbb{I}(w=0)/2 - \mathbb{I}(w=1)/2} \exp \left\{ -\frac{vz}{2\mu^2} - \frac{v}{2z} + \frac{v}{\mu} + \frac{1}{2} \log(v) - \mathbb{I}(w=1) \log(\mu) \right\}
 \end{aligned}$$

The natural parameters and sufficient statistics are $\{-\frac{v}{2\mu^2}, -\frac{v}{2}\}$ and $\{z, \frac{1}{z}\}$ respectively. \square

According to Balakrishnan & Kundu (2019), the expectation parameters are

$$\begin{aligned} m_1 &= \mathbb{E}_{q(z)}[z] = \mu + \frac{\mu^2}{2v} \\ m_2 &= \mathbb{E}_{q(z)}[z^{-1}] = \mu^{-1} + \frac{1}{2v} \end{aligned}$$

The sufficient statistics, natural parameters, and expectation parameters are summarized below:

$$\begin{bmatrix} z \\ \frac{1}{z} \end{bmatrix} \quad \begin{bmatrix} -\frac{v}{2\mu^2} \\ -\frac{v}{2} \end{bmatrix} \quad \begin{bmatrix} \mu + \frac{\mu^2}{2v} \\ \mu^{-1} + \frac{1}{2v} \end{bmatrix}$$

Lemma 7 *The joint distribution given in (40) is a minimal conditional-EF.*

Proof: Since λ_w is known in this case, we only need to show there exists an one-to-one mapping between the natural parameter and the expectation parameter. Just by observing the parameters given above, we can see that there exists an one-to-one mapping between the natural parameter and $\{\mu, v\}$ since $\mu > 0$ and $v > 0$. Furthermore, we know that $m_1 m_2 > 1$ and $m_1 > 0$. We can show that there also exists an one-to-one mapping between $\{\mu, v\}$ and the expectation parameter by noticing that

$$\begin{aligned} \mu &= \sqrt{m_1/m_2} \\ v &= \frac{1}{2(m_2 - \sqrt{m_2/m_1})} \end{aligned}$$

Since one-to-one mapping is transitive, we know that mapping between natural and expectation parameters is one-to-one. Hence proved. \square

Note that we can use the implicitly re-parametrization trick to compute the gradient w.r.t. μ and v . Furthermore, the expectation parameters m_1 and m_2 are functions of μ, v and the gradients can be obtained in terms of the gradient with respect to μ and v by using the chain rule.

D. Studnet's t-distribution

Lemma 8 *The joint distribution $\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})\mathcal{IG}(w|a, a)$, where $\mathbf{z} \in \mathcal{R}^d$, is a curved exponential family distribution.*

Proof: The joint-distribution can be expressed as a four-parameter exponential form as shown below:

$$\begin{aligned} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})\mathcal{IG}(w|a, a) &= \det(2\pi w\boldsymbol{\Sigma})^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T (w\boldsymbol{\Sigma})^{-1} (\mathbf{z} - \boldsymbol{\mu})\right\} \frac{a^a}{\Gamma(a)} w^{-a-1} \exp\left\{-\frac{a}{w}\right\} \\ &= (2\pi w)^{-d/2} w^{-1} \exp\left\{-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T (w\boldsymbol{\Sigma})^{-1} (\mathbf{z} - \boldsymbol{\mu}) - \frac{1}{2} \log \det(\boldsymbol{\Sigma})\right. \\ &\quad \left. - \frac{a}{w} - a \log w - (\log \Gamma(a) - a \log a)\right\} \\ &= (2\pi w)^{-d/2} w^{-1} \exp\left\{\left\langle -\frac{1}{2}\boldsymbol{\Sigma}^{-1}, w\mathbf{z}\mathbf{z}^T \right\rangle + \left\langle \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}, w\mathbf{z} \right\rangle + \left\langle -\frac{1}{2}\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}, w^{-1} \right\rangle\right. \\ &\quad \left. + \left\langle -a, w^{-1} + \log w \right\rangle - \left[\log \Gamma(a) - a \log a + \frac{1}{2} \log \det(\boldsymbol{\Sigma})\right]\right\} \\ &= (2\pi w)^{-d/2} w^{-1} \exp\left\{\left\langle \boldsymbol{\Lambda}_1, w\mathbf{z}\mathbf{z}^T \right\rangle + \left\langle \boldsymbol{\Lambda}_2, w\mathbf{z} \right\rangle + \left\langle \lambda_3, w^{-1} \right\rangle\right. \\ &\quad \left. + \left\langle \lambda_4, w^{-1} + \log w \right\rangle - \left[\log \Gamma(-\lambda_4) + \lambda_4 \log(-\lambda_4) - \frac{1}{2} \log \det(-2\boldsymbol{\Lambda}_1)\right]\right\}, \end{aligned}$$

where the following are the natural parameters:

$$\boldsymbol{\Lambda}_1 := -\frac{1}{2}\boldsymbol{\Sigma}^{-1}, \quad \boldsymbol{\Lambda}_2 := \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}, \quad \lambda_3 := -\frac{1}{2}\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}, \quad \lambda_4 := -a$$

We can see that λ_3 is fully determined by $\boldsymbol{\Lambda}_1$ and $\boldsymbol{\Lambda}_2$, i.e.,

$$\lambda_3 = -\frac{1}{2}\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} = -\frac{1}{2}(\boldsymbol{\Sigma}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu})^T (\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}) = -\frac{1}{2}\left((-2\boldsymbol{\Lambda}_1)^{-1}\boldsymbol{\Lambda}_2\right)^T \boldsymbol{\Lambda}_2$$

In a minimal representation we can specify all four parameters freely, but in this case we have less degree of freedom. Therefore, this is a curved EF representation. \square

Instead of using the above 4 parameter form, we can write the distribution in the conditional EF form given in Definition 2.

Lemma 9 *The joint distribution of Studnet's t-distribution given in (17) can be written in a conditional EF form.*

Proof: We can rewrite the conditional $q(\mathbf{z}|w)$ in an EF-form as follows:

$$\begin{aligned} q(\mathbf{z}|w) &= \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma}) \\ &= (2\pi)^{-d/2} \exp(\{\text{Tr}(-\frac{1}{2}\boldsymbol{\Sigma}^{-1}w^{-1}\mathbf{z}\mathbf{z}^T) + \boldsymbol{\mu}^T\boldsymbol{\Sigma}^{-1}w^{-1}\mathbf{z} - \frac{1}{2}(w^{-1}\boldsymbol{\mu}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \log \det(w\boldsymbol{\Sigma}))\}) \end{aligned}$$

The sufficient statistics $\phi_z(\mathbf{z}, w) = \{w^{-1}\mathbf{z}, w^{-1}\mathbf{z}\mathbf{z}^T\}$. The natural parameter is $\boldsymbol{\lambda}_z = \{\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}, -\frac{1}{2}\boldsymbol{\Sigma}^{-1}\}$. Since $q(w)$ is a inverse Gamma distribution, which is a EF distribution as shown below, the joint $q(\mathbf{z}, w)$ is a conditional EF. The EF form of the inverse gamma distribution is shown below:

$$q(w|a) = w^{-1} \exp(-(\log(w) + \frac{1}{w})a - (\log \Gamma(a) - a \log(a)))$$

We can read the sufficient statistics $\phi_w(w) = -\log(w) - \frac{1}{w}$ and the natural parameter $\lambda_w = a$ from this form. \square

Using the fact that $\mathbb{E}_{q(w|a)}[1/w] = a/a = 1$, and $\mathbb{E}_{q(w|a)}[\log w] = \log a - \psi(a)$, we can derive the expectation parameter shown below:

$$\begin{aligned} \mathbf{m} &:= \mathbb{E}_{q(w|a)q(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})}[w^{-1}\mathbf{z}] = \boldsymbol{\mu}, \\ \mathbf{M} &:= \mathbb{E}_{q(w|a)q(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})}[w^{-1}\mathbf{z}\mathbf{z}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma}, \\ m_a &:= \mathbb{E}_{q(w|a)}\left[-\frac{1}{w} - \log(w)\right] = -1 - \log a + \psi(a) \end{aligned}$$

The sufficient statistics, natural parameters, and expectation parameters are summarized below:

$$\begin{bmatrix} -1/w - \log w \\ \mathbf{z}/w \\ \mathbf{z}\mathbf{z}^T/w \end{bmatrix} \quad \begin{bmatrix} a \\ \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ -\frac{1}{2}\boldsymbol{\Sigma}^{-1} \end{bmatrix} \quad \begin{bmatrix} -1 - \log a + \psi(a) \\ \boldsymbol{\mu} \\ \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma} \end{bmatrix}$$

The following lemma shows that the Student's t-distribution is an MCEF obtained by establishing one-to-one mapping between natural and expectation parameters.

Lemma 10 *The joint distribution of Student's t-distribution shown in (17) is a minimal conditional-EF.*

Proof: The proof is rather simple. First we note that the expectation parameters for $q(\mathbf{z}|w)$ do not depend on $\lambda_w := a$. In fact, the mapping between the last two natural and expectation parameter is one-to-one since they correspond to a Gaussian distribution which has a minimal representation.

The only thing remaining is to show that the mapping between a and $m_w(a) := -1 - \log a + \psi(a)$ is one-to-one. Since $\nabla_a m_w(a)$ is the Fisher information of $q(w)$, we can show this when $\nabla_a m_w(a) > 0$. The gradient $\nabla_a m_w(a) = \nabla_a \psi(a) - 1/a$. According to Eq. 1.4 in Batir (2005), we have $\nabla_a \psi(a) - 1/a - 1/(2a^2) > 0$ when $a > 0$. Therefore, $\nabla_a \psi(a) - 1/a > 0$ which establishes that the Fisher information is positive, therefore the distribution is a minimal EF. This completes the proof. \square

D.1. Derivation of the NGD Update

Let's consider $q(w) = \mathcal{IG}(w|a, a)$ and $q(\mathbf{z}|w) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})$. We denote the log-likelihood for the n 'th data point by $f_n(\mathbf{z}) := -\log p(\mathcal{D}_n|\mathbf{z})$ with a Student's t-prior on \mathbf{z} expressed as a scale mixture of Gaussians $p(\mathbf{z}, w) =$

$\mathcal{IG}(w|a_0, a_0)\mathcal{N}(\mathbf{z}|\mathbf{0}, w\mathbf{I})$. We use the lower bound defined in the joint-distribution $p(\mathcal{D}, \mathbf{z}, w)$:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\lambda}) &= \mathbb{E}_{q(\mathbf{z}, w)} [\log p(\mathcal{D}, \mathbf{z}, w) - \log q(\mathcal{D}, \mathbf{z}, w)] \\ &= \mathbb{E}_{q(\mathbf{z}, w)} \left[\sum_{n=1}^N \underbrace{\log p(\mathcal{D}_n|\mathbf{z})}_{:= -f_n(\mathbf{z})} + \log \frac{\mathcal{N}(\mathbf{z}|\mathbf{0}, w\mathbf{I})}{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})} + \log \frac{\mathcal{IG}(w|a_0, a_0)}{\mathcal{IG}(w|a, a)} \right],\end{aligned}$$

Our goal is to compute the gradient of this ELBO with respect to the expectation parameters.

Since the expectation parameters \mathbf{m}_z only depend on $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, we can write the gradient with respect to them using the chain rule (similar to the finite mixture of Gaussians case):

$$\begin{aligned}\nabla_m \mathcal{L}(\boldsymbol{\lambda}) &= \nabla_{\boldsymbol{\mu}} \mathcal{L}(\boldsymbol{\lambda}) - 2\nabla_{\boldsymbol{\Sigma}} \mathcal{L}(\boldsymbol{\lambda})\boldsymbol{\mu} \\ \nabla_M \mathcal{L}(\boldsymbol{\lambda}) &= \nabla_{\boldsymbol{\Sigma}} \mathcal{L}(\boldsymbol{\lambda})\end{aligned}$$

These gradients of the lower bound can be obtained as follows:

$$\begin{aligned}\nabla_{\boldsymbol{\mu}} \mathcal{L}(\boldsymbol{\lambda}) &= - \sum_{n=1}^N \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(\mathbf{z}, w)} [f_n(\mathbf{z})] - \boldsymbol{\mu} \\ \nabla_{\boldsymbol{\Sigma}} \mathcal{L}(\boldsymbol{\lambda}) &= - \sum_{n=1}^N \nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{q(\mathbf{z}, w)} [f_n(\mathbf{z})] - \frac{1}{2}\mathbf{I} + \frac{1}{2}\boldsymbol{\Sigma}^{-1}\end{aligned}$$

Plugging these in the natural-gradient updates (14), we get the following updates (we have simplified these in the same way as explained in Appendix B.2; more details in Khan et al. (2018)):

$$\begin{aligned}\boldsymbol{\Sigma}^{-1} &\leftarrow (1 - \beta)\boldsymbol{\Sigma}^{-1} + 2\beta \sum_{n=1}^N \nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{q(\mathbf{z}, w)} [f_n(\mathbf{z})] + \beta\mathbf{I} \\ \boldsymbol{\mu} &\leftarrow \boldsymbol{\mu} - \beta \sum_{n=1}^N \boldsymbol{\Sigma} (\nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(\mathbf{z})} [f_n(\mathbf{z})] + \boldsymbol{\mu})\end{aligned}$$

To compute the gradients, the reparametrization trick (Kingma & Welling, 2013) can be used. However, we can do better by using the extended Bonnet's and Price's theorems for Student's t-distribution (Lin et al., 2019). Assuming that $f_n(\mathbf{z})$ satisfies the assumptions needed for these two theorems, we can use the following stochastic-gradient approximations for the gradients:

$$\begin{aligned}\nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(w)\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})} [f(\mathbf{z})] &= \mathbb{E}_{q(\mathbf{z})} [\nabla_{\mathbf{z}} f(\mathbf{z})] \approx \nabla_{\mathbf{z}} f(\mathbf{z}), \\ \nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{q(w)\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})} [f(\mathbf{z})] &= \frac{1}{2} \mathbb{E}_{q(\mathbf{z})} [u(\mathbf{z}) \nabla_{\mathbf{z}}^2 f(\mathbf{z})] \approx \frac{1}{2} u(\mathbf{z}) \nabla_{\mathbf{z}}^2 f(\mathbf{z}), \\ &= \frac{1}{2} \mathbb{E}_{q(w, \mathbf{z})} [w \nabla_{\mathbf{z}}^2 f(\mathbf{z})] \approx \frac{1}{2} w \nabla_{\mathbf{z}}^2 f(\mathbf{z}),\end{aligned}$$

where $\mathbf{z} \in \mathcal{R}^d$ is generated from $q(\mathbf{z})$, w is generated from $q(w)$, and

$$u(\mathbf{z}) := \frac{a + \frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{z} - \boldsymbol{\mu})}{(a + d/2 - 1)}.$$

Using these gradient, we get the following update:

$$\begin{aligned}\boldsymbol{\Sigma}^{-1} &\leftarrow (1 - \beta)\boldsymbol{\Sigma}^{-1} + \beta [u(\mathbf{z}) \nabla_{\mathbf{z}}^2 f_n(\mathbf{z}) + \mathbf{I}/N], \\ \boldsymbol{\mu} &\leftarrow \boldsymbol{\mu} - \beta \boldsymbol{\Sigma} [\nabla_{\mathbf{z}} f_n(\mathbf{z}) + \boldsymbol{\mu}/N].\end{aligned}$$

Now we derive the NGD update for a . Recall that the natural parameter is a and the expectation parameter is $m_a = -1 - \log a + \psi(a)$. The gradient of the lower bound can be expressed as

$$\nabla_{m_a} \mathcal{L}(\boldsymbol{\lambda}) = - \sum_{n=1}^N \nabla_{m_a} \mathbb{E}_{q(\mathbf{z}, w)} [f_n(\mathbf{z})] + a_0 - a$$

which gives us the following update:

$$a \leftarrow (1 - \beta)a + \beta \left(a_0 - \sum_{n=1}^N \nabla_{m_a} \mathbb{E}_{q(z,w)} [f_n(\mathbf{z})] \right) \quad (41)$$

While the gradient with respect to the expectation parameter does not admit a closed-form expression, we can compute the gradient using the re-parametrization trick. According to (39), the gradient $\nabla_{m_a} \mathbb{E}_{q(z,w)} [f_n(\mathbf{z})]$ can be computed as

$$\nabla_{m_a} \mathbb{E}_{q(z,w)} [f_n(\mathbf{z})] = (\nabla_a m_a)^{-1} \nabla_a \mathbb{E}_{q(z,w)} [f_n(\mathbf{z})] = (\nabla_a \mathbb{E}_{q(w)} [\phi_w(w)])^{-1} \nabla_a \mathbb{E}_{q(z,w)} [f_n(\mathbf{z})]$$

Note that $\nabla_a \mathbb{E}_{q(w)} [\phi_w(w)] = \nabla_a (m_a)$ has a closed-form expression. However we have found that using stochastic approximation for both the numerator and denominator works better. Salimans & Knowles (2013) show that such approximation reduces the variance but introduce a bit bias. We use the reparameterization trick for both terms. Since $q(w)$ is (implicitly) re-parameterizable (Salimans & Knowles, 2013; Figurnov et al., 2018), the gradient can be computed as

$$\nabla_a \mathbb{E}_{q(w)} [\phi_w(w)] = - \int \mathcal{IG}(w|a, a) (\nabla_w [w^{-1} + \log w]) (\nabla_a w) dw \approx (w^{-2} - w^{-1}) \nabla_a w$$

where w is generated from $\mathcal{IG}(w|a, a)$. Similarly,

$$\begin{aligned} \nabla_a \mathbb{E}_{q(z,w)} [f_n(\mathbf{z})] &= \int \int q(w) (\nabla_w q(\mathbf{z}|w)) (\nabla_a w) f_n(\mathbf{z}) dw d\mathbf{z} \\ &= \int \int \mathcal{IG}(w|a, a) (\nabla_w \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})) (\nabla_a w) f_n(\mathbf{z}) dw d\mathbf{z} \end{aligned}$$

For stochastic approximation, we generate w from $q(w)$ and let $\hat{\boldsymbol{\Sigma}} = w\boldsymbol{\Sigma}$. The above expression can be approximated as below.

$$\begin{aligned} \nabla_a \mathbb{E}_{q(z,w)} [f_n(\mathbf{z})] &\approx \int (\nabla_w \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w\boldsymbol{\Sigma})) (\nabla_a w) f_n(\mathbf{z}) d\mathbf{z} \\ &= \int \text{Tr} \left(\boldsymbol{\Sigma} \nabla_{\hat{\boldsymbol{\Sigma}}} \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \hat{\boldsymbol{\Sigma}}) \right) (\nabla_a w) f_n(\mathbf{z}) d\mathbf{z} \\ &= \nabla_a w \text{Tr} \left(\boldsymbol{\Sigma} \nabla_{\hat{\boldsymbol{\Sigma}}} \mathbb{E}_{\mathcal{N}(z|\boldsymbol{\mu}, \hat{\boldsymbol{\Sigma}})} [f_n(\mathbf{z})] \right) \\ &= \frac{\nabla_a w}{2} \text{Tr} \left(\boldsymbol{\Sigma} \mathbb{E}_{\mathcal{N}(z|\boldsymbol{\mu}, \hat{\boldsymbol{\Sigma}})} [\nabla_z^2 f_n(\mathbf{z})] \right), \end{aligned}$$

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

We then use stochastic approximation to get the desired update.

$$\nabla_a \mathbb{E}_{q(z,w)} [f_n(\mathbf{z})] \approx \frac{\nabla_a w}{2} \text{Tr} \left(\boldsymbol{\Sigma} \nabla_z^2 f_n(\mathbf{z}) \right)$$

where \mathbf{z} is generated from $q(\mathbf{z})$ and w is generated from $q(w)$.

Another example is the symmetric normal inverse-Gaussian distribution, which can be found at Appendix H.

E. Multivariate Skew-Gaussian Distribution

We consider the following variational distribution .

$$q(\mathbf{z}, w) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu} + |w|\boldsymbol{\alpha}, \boldsymbol{\Sigma}) \mathcal{N}(w|0, 1) \quad (42)$$

The marginal distribution is known as the multivariate skew Gaussian distribution (Azzalini, 2005) as shown in the following lemma.

Lemma 11 *The marginal distribution $q(\mathbf{z})$ is $2\Phi \left(\frac{(\mathbf{z}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}{\sqrt{1+\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}} \right) \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma} + \boldsymbol{\alpha} \boldsymbol{\alpha}^T)$, where $\Phi(\cdot)$ is CDF of the standard univariate Gaussian distribution.*

Proof: The marginal distribution \mathbf{z} is

$$\begin{aligned} q(\mathbf{z}) &= 2 \int_0^{+\infty} \mathcal{N}(w|0, 1) \mathcal{N}(\mathbf{z}|\boldsymbol{\mu} + w\boldsymbol{\alpha}, \boldsymbol{\Sigma}) dw \\ &= 2 \int_0^{+\infty} \mathcal{N}(w|0, 1) \mathcal{N}(\mathbf{z} - w\boldsymbol{\alpha}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) dw \end{aligned}$$

By grouping terms related to w together and completing a Gaussian form for w , we obtain the following expression

$$\begin{aligned} q(\mathbf{z}) &= 2 \int_0^{+\infty} \mathcal{N}(w|\frac{(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}{1 + \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}, (1 + \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha})^{-1}) \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma} + \boldsymbol{\alpha} \boldsymbol{\alpha}^T) dw \\ &= 2\Phi\left(\frac{(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}{\sqrt{1 + \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}}\right) \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma} + \boldsymbol{\alpha} \boldsymbol{\alpha}^T) \end{aligned}$$

where we move from the first step to the second step using the fact that

$$\begin{aligned} \int_0^{+\infty} \mathcal{N}(w|u, \sigma^2) dw &= \int_{-u/\sigma}^{+\infty} \mathcal{N}(w'|0, 1) dw' \\ &= 1 - \Phi(-u/\sigma) \\ &= \Phi(u/\sigma). \end{aligned}$$

□

Lemma 12 *The joint distribution of skew-Gaussian distribution given at Eq. (42) can be written in a conditional EF form.*

Proof: We first rewrite $q(\mathbf{z}|w)$ in a EF-form as follows:

$$\begin{aligned} q(\mathbf{z}|w) &= \mathcal{N}(\mathbf{z}|\boldsymbol{\mu} + |w|\boldsymbol{\alpha}, \boldsymbol{\Sigma}) \\ &= (2\pi)^{-d/2} \exp\left(\left\{\text{Tr}\left(-\frac{1}{2}\boldsymbol{\Sigma}^{-1}\mathbf{z}\mathbf{z}^T\right) + |w|\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \mathbf{z} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{z} - \frac{1}{2}((\boldsymbol{\mu} + |w|\boldsymbol{\alpha})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} + |w|\boldsymbol{\alpha}) + \log \det(\boldsymbol{\Sigma}))\right\}\right) \end{aligned}$$

The sufficient statistics $\phi_z(\mathbf{z}, w) = \{\mathbf{z}, |w|\mathbf{z}, \mathbf{z}\mathbf{z}^T\}$ and the natural parameter $\boldsymbol{\lambda}_z = \{\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}\boldsymbol{\alpha}, -\frac{1}{2}\boldsymbol{\Sigma}^{-1}\}$ can be read from the form. Since $q(w)$ is a univariate Gaussian distribution with known parameters, which is a EF distribution. Therefore, the joint distribution $q(\mathbf{z}, w)$ is a conditional EF. □

Let $c = \sqrt{\frac{2}{\pi}}$. Using the fact that $\mathbb{E}_{q(w|0,1)}[|w|] = c$, we can derive the expectation parameter shown below:

$$\begin{aligned} \mathbf{m} &:= \mathbb{E}_{\mathcal{N}(w|0,1), \mathcal{N}(z|\boldsymbol{\mu}+|w|\boldsymbol{\alpha}, \boldsymbol{\Sigma})} [\mathbf{z}] = \boldsymbol{\mu} + c\boldsymbol{\alpha}, \\ \mathbf{m}_\alpha &:= \mathbb{E}_{\mathcal{N}(w|0,1), \mathcal{N}(z|\boldsymbol{\mu}+|w|\boldsymbol{\alpha}, \boldsymbol{\Sigma})} [|w|\mathbf{z}] = c\boldsymbol{\mu} + \boldsymbol{\alpha}, \\ \mathbf{M} &:= \mathbb{E}_{\mathcal{N}(w|0,1), \mathcal{N}(z|\boldsymbol{\mu}+|w|\boldsymbol{\alpha}, \boldsymbol{\Sigma})} [\mathbf{z}\mathbf{z}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\alpha}\boldsymbol{\alpha}^T + c(\boldsymbol{\mu}\boldsymbol{\alpha}^T + \boldsymbol{\alpha}\boldsymbol{\mu}^T) + \boldsymbol{\Sigma} \end{aligned}$$

The sufficient statistics, natural parameters, and expectation parameters are summarized below:

$$\begin{bmatrix} \mathbf{z} \\ |w|\mathbf{z} \\ \mathbf{z}\mathbf{z}^T \end{bmatrix} \quad \begin{bmatrix} \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ \boldsymbol{\Sigma}^{-1}\boldsymbol{\alpha} \\ -\frac{1}{2}\boldsymbol{\Sigma}^{-1} \end{bmatrix} \quad \begin{bmatrix} \boldsymbol{\mu} + c\boldsymbol{\alpha} \\ c\boldsymbol{\mu} + \boldsymbol{\alpha} \\ \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\alpha}\boldsymbol{\alpha}^T + c(\boldsymbol{\mu}\boldsymbol{\alpha}^T + \boldsymbol{\alpha}\boldsymbol{\mu}^T) + \boldsymbol{\Sigma} \end{bmatrix}$$

The following lemma shows that the skew-Gaussian distribution is indeed a minimal conditional-EF.

Lemma 13 *Multivariate skew Gaussians is a minimal conditional-EF.*

Proof: Since λ_w is known in this case, we only need to show there exists an one-to-one mapping between the natural parameter and the expectation parameter. Just by observing the parameters given above, we can see that there exists an one-to-one mapping between the natural parameter and $\{\mu, \alpha, \Sigma\}$. We can show that there also exists an one-to-one mapping between $\{\mu, \alpha, \Sigma\}$ and the expectation parameter by noticing that

$$\mu = \frac{\mathbf{m} - c\mathbf{m}_\alpha}{1 - c^2} \quad (43)$$

$$\alpha = \frac{\mathbf{m}_\alpha - c\mathbf{m}}{1 - c^2} \quad (44)$$

$$\Sigma = \mathbf{M} - \frac{\mathbf{m}\mathbf{m}^T + \mathbf{m}_\alpha\mathbf{m}_\alpha^T - c(\mathbf{m}_\alpha\mathbf{m}^T + \mathbf{m}\mathbf{m}_\alpha^T)}{1 - c^2} \quad (45)$$

Since one-to-one mapping is transitive, we know that mapping between natural and expectation parameters is one-to-one. Hence proved. \square

E.1. Derivation of the NGD Update

Let's consider the variational approximation using the skew-Gaussian distribution $q(\mathbf{z}|\lambda)$. We consider the following model with a Gaussian-prior $\mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})$ on \mathbf{z} , where the log-likelihood for the n 'th data point is denoted by $p(\mathcal{D}_n|\mathbf{z})$.

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

We use the lower bound defined in the following distribution $p(\mathcal{D}, \mathbf{z})$:

$$\mathcal{L}(\lambda) = \mathbb{E}_{q(\mathbf{z}|\lambda)} \left[\sum_{n=1}^N \underbrace{\log p(\mathcal{D}_n|\mathbf{z})}_{:= -f_n(\mathbf{z})} + \log \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I}) - \log q(\mathbf{z}|\lambda) \right],$$

where $q(\mathbf{z}) = 2\Phi\left(\frac{(\mathbf{z}-\mu)^T \Sigma^{-1} \alpha}{\sqrt{1+\alpha^T \Sigma^{-1} \alpha}}\right) \mathcal{N}(\mathbf{z}|\mu, \Sigma + \alpha\alpha^T)$ and recall that $\Phi(\cdot)$ denotes the CDF of the standard univariate normal distribution.

Our goal is to compute the gradient of this ELBO with respect to the expectation parameters.

E.2. Natural Gradient for $q(\mathbf{z}|w)$

We do not need to compute these gradients with respect to the expectation parameters directly. The gradients can be computed in terms of $\{\mu, \alpha, \Sigma\}$.

Using the mapping at (43)-(45) and the chain rule, we can express the following gradients with respect to the expectation parameters in terms of the gradients with respect to μ, α , and Σ .

$$\begin{aligned} \nabla_m \mathcal{L} &= \frac{1}{1-c^2} \nabla_\mu \mathcal{L} - \frac{c}{1-c^2} \nabla_\alpha \mathcal{L} - 2(\nabla_\Sigma \mathcal{L}) \mu \\ \nabla_{m_\alpha} \mathcal{L} &= \frac{1}{1-c^2} \nabla_\alpha \mathcal{L} - \frac{c}{1-c^2} \nabla_\mu \mathcal{L} - 2(\nabla_\Sigma \mathcal{L}) \alpha \\ \nabla_M \mathcal{L} &= \nabla_\Sigma \mathcal{L} \end{aligned}$$

By plugging the gradients into the update in (14) and then re-expressing the update in terms of μ, α, Σ (we have simplified these in the same way as explained in Appendix B.2), we obtain the natural gradient update in terms of μ, α , and Σ .

$$\Sigma^{-1} \leftarrow \Sigma^{-1} - 2\beta \nabla_\Sigma \mathcal{L} \quad (46)$$

$$\mu \leftarrow \mu + \beta \Sigma \left(\frac{1}{1-c^2} \nabla_\mu \mathcal{L} - \frac{c}{1-c^2} \nabla_\alpha \mathcal{L} \right) \quad (47)$$

$$\alpha \leftarrow \alpha + \beta \Sigma \left(\frac{1}{1-c^2} \nabla_\alpha \mathcal{L} - \frac{c}{1-c^2} \nabla_\mu \mathcal{L} \right) \quad (48)$$

Recall that the lower bound is

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})} \left[- \sum_{n=1}^N f_n(\mathbf{z}) + \underbrace{\log \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})}_{\text{prior}} - \underbrace{\log q(\mathbf{z}|\boldsymbol{\lambda})}_{\text{entropy}} \right],$$

For the prior term, there is a closed-form expression for gradient computation.

$$\mathbb{E}_{q(\mathbf{z})} [\log \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})] = -\frac{d}{2} \log(2\pi) + \frac{d\delta}{2} - \frac{\delta}{2} (\boldsymbol{\alpha}^T \boldsymbol{\alpha} + 2c\boldsymbol{\mu}^T \boldsymbol{\alpha} + \text{Tr}(\boldsymbol{\Sigma}) + \boldsymbol{\mu}^T \boldsymbol{\mu})$$

It is easy to show that the gradients about the prior term are

$$\begin{aligned} \mathbf{g}_{\boldsymbol{\mu}}^{\text{prior}} &= -\delta (\boldsymbol{\mu} + c\boldsymbol{\alpha}) \\ \mathbf{g}_{\boldsymbol{\alpha}}^{\text{prior}} &= -\delta (\boldsymbol{\alpha} + c\boldsymbol{\mu}) \\ \mathbf{g}_{\boldsymbol{\Sigma}}^{\text{prior}} &= -\frac{\delta}{2} \mathbf{I} \end{aligned}$$

For the entropy term, by Contreras-Reyes & Arellano-Valle (2012); Arellano-Valle et al. (2013), it can be expressed as follows.

$$\mathbb{E}_{q(\mathbf{z})} [-\log q(\mathbf{z}|\boldsymbol{\lambda})] = \frac{d}{2} (\log(2\pi) + 1) + \frac{1}{2} \log |\boldsymbol{\Sigma} + \boldsymbol{\alpha}\boldsymbol{\alpha}^T| - 2\mathbb{E}_{\mathcal{N}(z_3|0, \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha})} [\Phi(z_3) \log(\Phi(z_3))] - \log(2) \quad (49)$$

We can use the re-parametrization trick to compute the gradients about the entropy term. However, we can found out the exact gradients usually works better. Using the expression in (49), the gradients of the entropy term are given as follows:

$$\mathbf{g}_{\boldsymbol{\mu}}^{\text{entropy}} = 0 \quad (50)$$

$$\mathbf{g}_{\boldsymbol{\alpha}}^{\text{entropy}} = (\boldsymbol{\Sigma} + \boldsymbol{\alpha}\boldsymbol{\alpha}^T)^{-1} \boldsymbol{\alpha} - \mathbb{E}_{\mathcal{N}(z_3|0, \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} / (1 + \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}))} \left[\frac{\log(\Phi(z_3))}{\sqrt{2\pi(1 + \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha})}} \frac{2z_3 \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}} \right] \quad (51)$$

$$\mathbf{g}_{\boldsymbol{\Sigma}}^{\text{entropy}} = \frac{1}{2} (\boldsymbol{\Sigma} + \boldsymbol{\alpha}\boldsymbol{\alpha}^T)^{-1} + \mathbb{E}_{\mathcal{N}(z_3|0, \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} / (1 + \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}))} \left[\frac{\log(\Phi(z_3))}{\sqrt{2\pi(1 + \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha})}} \frac{z_3 \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1}}{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}} \right] \quad (52)$$

where the expectations involve 1d integrations, which can be computed by Gauss-Hermite quadrature.

The remaining thing is to compute the gradients about $\mathbb{E}_{q(\mathbf{z})} [f_n(\mathbf{z})]$. To compute the gradients, the reparametrization trick can be used. However, we can do better by using the extended Bonnet's and Price's theorems for skew-Gaussian distribution (Lin et al., 2019). Assuming that $f_n(\mathbf{z})$ satisfies the assumptions needed for these two theorems, we obtain the following gradient expression:

$$\begin{aligned} \mathbf{g}_1^n &:= \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(\mathbf{z})} [f_n(\mathbf{z})] = \mathbb{E}_{q(\mathbf{z})} [\nabla_{\mathbf{z}} f_n(\mathbf{z})] \approx \nabla_{\mathbf{z}} f_n(\mathbf{z}) \\ \mathbf{g}_2^n &:= \nabla_{\boldsymbol{\alpha}} \mathbb{E}_{q(\mathbf{z})} [f_n(\mathbf{z})] \\ &= \mathbb{E}_{q(\mathbf{z})} [u(\mathbf{z}) \nabla_{\mathbf{z}} f_n(\mathbf{z})] + v \mathbb{E}_{\mathcal{N}(\hat{\mathbf{z}}|\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\nabla_{\hat{\mathbf{z}}} f_n(\hat{\mathbf{z}})] \approx u(\mathbf{z}) \nabla_{\mathbf{z}} f_n(\mathbf{z}) + v \nabla_{\hat{\mathbf{z}}} f_n(\hat{\mathbf{z}}) \\ &= \mathbb{E}_{q(w, \mathbf{z})} [|w| \nabla_{\mathbf{z}} f_n(\mathbf{z})] \approx |w| \nabla_{\mathbf{z}} f_n(\mathbf{z}) \\ \mathbf{g}_3^n &:= 2\nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{q(\mathbf{z})} [f_n(\mathbf{z})] = \mathbb{E}_{q(\mathbf{z})} [\nabla_{\mathbf{z}}^2 f_n(\mathbf{z})] \approx \nabla_{\mathbf{z}}^2 f_n(\mathbf{z}) \end{aligned}$$

where $v = \frac{c}{(1 + \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha})}$, $u(\mathbf{z}) := \frac{(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}{1 + \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}$, and $w \sim \mathcal{N}(w|0, 1)$, $\hat{\mathbf{z}} \sim \mathcal{N}(\hat{\mathbf{z}}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$, $\mathbf{z} = \hat{\mathbf{z}} + |w|\boldsymbol{\alpha}$.

Putting together, we can express the gradients in the following form, which will be used for deriving the extended variational

Adam update.

$$\begin{aligned}\nabla_{\mu}\mathcal{L}(\lambda) &= -\sum_{n=1}^N \left(\underbrace{\mathbf{g}_1^n - \mathbf{g}_{\mu}^{\text{entropy}}/N}_0 \right) \underbrace{-\delta(\mu + c\alpha)}_{\mathbf{g}_{\mu}^{\text{prior}}} \\ &= -\sum_{n=1}^N \underbrace{\mathbf{g}_1^n}_{:=\mathbf{g}_{\mu}^n} - \delta(\mu + c\alpha)\end{aligned}\quad (53)$$

$$\nabla_{\alpha}\mathcal{L}(\lambda) = -\sum_{n=1}^N \left(\underbrace{\mathbf{g}_2^n - \mathbf{g}_{\alpha}^{\text{entropy}}/N}_{:=\mathbf{g}_{\alpha}^n} \right) \underbrace{-\delta(\alpha + c\mu)}_{\mathbf{g}_{\alpha}^{\text{prior}}} \quad (54)$$

$$\begin{aligned}\nabla_{\Sigma}\mathcal{L}(\lambda) &= -\frac{1}{2} \sum_{n=1}^N (\mathbf{g}_3^n - 2\mathbf{g}_{\Sigma}^{\text{entropy}}/N) \underbrace{-\frac{\delta}{2}\mathbf{I} + \frac{1}{2}(\Sigma^{-1} - \Sigma^{-1})}_{\mathbf{g}_{\Sigma}^{\text{prior}}} \\ &= -\frac{1}{2} \sum_{n=1}^N \left(\underbrace{\mathbf{g}_3^n - 2\mathbf{g}_{\Sigma}^{\text{entropy}}/N + \Sigma^{-1}/N}_{:=\mathbf{g}_{\Sigma}^n} \right) - \frac{\delta}{2}\mathbf{I} + \frac{1}{2}\Sigma^{-1}\end{aligned}\quad (55)$$

For stochastic approximation, we can sub-sampling a data point n and use MC samples to approximate \mathbf{g}_1^n , \mathbf{g}_2^n , and \mathbf{g}_3^n . Plugging these stochastic gradients into (46)-(48), we obtain the NGD update:

$$\begin{aligned}\Sigma^{-1} &\leftarrow (1 - \beta)\Sigma^{-1} + \beta(\delta\mathbf{I} + N\mathbf{g}_{\Sigma}^n) \\ \mu &\leftarrow \mu - \beta\Sigma\left(\frac{N}{1 - c^2}(\mathbf{g}_{\mu}^n - c\mathbf{g}_{\alpha}^n) + \delta\mu\right) \\ \alpha &\leftarrow \alpha - \beta\Sigma\left(\frac{N}{1 - c^2}(\mathbf{g}_{\alpha}^n - c\mathbf{g}_{\mu}^n) + \delta\alpha\right)\end{aligned}$$

F. Multivariate Exponentially Modified Gaussian Distribution

We consider the following mixture distribution.

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

The marginal distribution is a multivariate extension of the exponentially modified Gaussian distribution (Grushka, 1972) and the Gaussian minus exponential distribution (Carr & Madan, 2009) due to the following lemma.

Lemma 14 *The marginal distribution $q(\mathbf{z})$ ($\alpha \neq \mathbf{0}$) is*

$$q(\mathbf{z}) = \frac{\sqrt{2\pi}\det(2\pi\Sigma)^{-\frac{1}{2}}}{\sqrt{\alpha^T\Sigma^{-1}\alpha}}\Phi\left(\frac{(\mathbf{z} - \mu)^T\Sigma^{-1}\alpha - 1}{\sqrt{\alpha^T\Sigma^{-1}\alpha}}\right)\exp\left\{\frac{1}{2}\left[\frac{\left((\mathbf{z} - \mu)^T\Sigma^{-1}\alpha - 1\right)^2}{\alpha^T\Sigma^{-1}\alpha} - (\mathbf{z} - \mu)^T\Sigma^{-1}(\mathbf{z} - \mu)\right]\right\},$$

where $\Phi(\cdot)$ is the CDF of the standard univariate Gaussian distribution.

In the univariate case, the marginal distribution becomes the exponentially modified Gaussian distribution when $\alpha > 0$ and the Gaussian minus exponential distribution when $\alpha < 0$.

Proof: The marginal distribution \mathbf{z} is

$$q(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\alpha}, \boldsymbol{\Sigma}) = \int_0^{+\infty} \text{Exp}(w|0, 1) \mathcal{N}(\mathbf{z}|\boldsymbol{\mu} + w\boldsymbol{\alpha}, \boldsymbol{\Sigma}) dw$$

By grouping terms related to w together and completing a Gaussian form for w , we obtain the following expression.

$$= \frac{\det(2\pi\boldsymbol{\Sigma})^{-\frac{1}{2}}}{\sqrt{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}} \int_0^{+\infty} \mathcal{N}\left(w \middle| \frac{(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} - 1}{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}, \frac{1}{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}\right) \exp \left\{ \frac{1}{2} \left[\frac{\left((\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} - 1 \right)^2}{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}} - (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu}) \right] \right\} dw \quad (56)$$

$$= \frac{\sqrt{2\pi} \det(2\pi\boldsymbol{\Sigma})^{-\frac{1}{2}}}{\sqrt{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}} \Phi \left(\frac{(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} - 1}{\sqrt{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}} \right) \exp \left\{ \frac{1}{2} \left[\frac{\left((\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} - 1 \right)^2}{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}} - (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu}) \right] \right\} \quad (57)$$

where we move from Eq. (56) to Eq. (57) using the fact that $\int_0^{+\infty} \mathcal{N}(w|u, \sigma^2) dw = \Phi(u/\sigma)$. \square

Similar to the skew-Gaussian case, in this example, the sufficient statistics $\phi_z(\mathbf{z}, w) = \{\mathbf{z}, w\mathbf{z}, \mathbf{z}\mathbf{z}^T\}$ and the natural parameter $\boldsymbol{\lambda}_z = \{\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}\boldsymbol{\alpha}, -\frac{1}{2}\boldsymbol{\Sigma}^{-1}\}$ can be read from $q(\mathbf{z}|w) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu} + w\boldsymbol{\alpha}, \boldsymbol{\Sigma})$. The joint distribution $q(\mathbf{z}, w)$ is a conditional EF since $q(w)$ is an exponential distribution with known parameters, which is an EF distribution. Likewise, we can show that the joint distribution is also a minimal conditional EF. We can derive the expectation parameter shown below:

$$\begin{aligned} \mathbf{m} &:= \mathbb{E}_{\text{Exp}(w|1)\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}+w\boldsymbol{\alpha}, \boldsymbol{\Sigma})} [\mathbf{z}] = \boldsymbol{\mu} + \boldsymbol{\alpha}, \\ \mathbf{m}_\alpha &:= \mathbb{E}_{\text{Exp}(w|1)\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}+w\boldsymbol{\alpha}, \boldsymbol{\Sigma})} [w\mathbf{z}] = \boldsymbol{\mu} + 2\boldsymbol{\alpha}, \\ \mathbf{M} &:= \mathbb{E}_{\text{Exp}(w|1)\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}+w\boldsymbol{\alpha}, \boldsymbol{\Sigma})} [\mathbf{z}\mathbf{z}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + 2\boldsymbol{\alpha}\boldsymbol{\alpha}^T + (\boldsymbol{\mu}\boldsymbol{\alpha}^T + \boldsymbol{\alpha}\boldsymbol{\mu}^T) + \boldsymbol{\Sigma} \end{aligned}$$

The sufficient statistics, natural parameters, and expectation parameters are summarized below:

$$\begin{bmatrix} \mathbf{z} \\ w\mathbf{z} \\ \mathbf{z}\mathbf{z}^T \end{bmatrix} \quad \begin{bmatrix} \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ \boldsymbol{\Sigma}^{-1}\boldsymbol{\alpha} \\ -\frac{1}{2}\boldsymbol{\Sigma}^{-1} \end{bmatrix} \quad \begin{bmatrix} \boldsymbol{\mu} + \boldsymbol{\alpha} \\ \boldsymbol{\mu} + 2\boldsymbol{\alpha} \\ \boldsymbol{\mu}\boldsymbol{\mu}^T + 2\boldsymbol{\alpha}\boldsymbol{\alpha}^T + (\boldsymbol{\mu}\boldsymbol{\alpha}^T + \boldsymbol{\alpha}\boldsymbol{\mu}^T) + \boldsymbol{\Sigma} \end{bmatrix}$$

F.1. Derivation of the NGD Update

As shown in Appendix E.1, we consider the variational approximation using the exponentially modified Gaussian distribution $q(\mathbf{z}|\boldsymbol{\lambda})$. We consider the same model as Appendix E.1 with a Gaussian-prior $\mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})$ on \mathbf{z} . The lower bound is defined as below:

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})} \left[\sum_{n=1}^N \underbrace{\log p(\mathcal{D}_n|\mathbf{z})}_{:= -f_n(\mathbf{z})} + \log \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I}) - \log q(\mathbf{z}|\boldsymbol{\lambda}) \right],$$

where $q(\mathbf{z}) = \frac{\sqrt{2\pi} \det(2\pi\boldsymbol{\Sigma})^{-\frac{1}{2}}}{\sqrt{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}} \Phi \left(\frac{(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} - 1}{\sqrt{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}} \right) \exp \left\{ \frac{1}{2} \left[\frac{\left((\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} - 1 \right)^2}{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}} - (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu}) \right] \right\}$ and recall that $\Phi(\cdot)$ denotes the CDF of the standard univariate normal distribution.

Our goal is to compute the gradient of this ELBO with respect to the expectation parameters.

F.2. Natural Gradient for $q(\mathbf{z}|w)$

We do not need to compute these gradients with respect to the expectation parameters directly. The gradients can be computed in terms of $\{\boldsymbol{\mu}, \boldsymbol{\alpha}, \boldsymbol{\Sigma}\}$.

Similarly, by the chain rule, we can express the following gradients with respect to the expectation parameters in terms of the gradients with respect to $\boldsymbol{\mu}$, $\boldsymbol{\alpha}$, and $\boldsymbol{\Sigma}$.

$$\begin{aligned}\nabla_m \mathcal{L} &= 2\nabla_\mu \mathcal{L} - \nabla_\alpha \mathcal{L} - 2(\nabla_\Sigma \mathcal{L}) \boldsymbol{\mu} \\ \nabla_{m_\alpha} \mathcal{L} &= \nabla_\alpha \mathcal{L} - \nabla_\mu \mathcal{L} - 2(\nabla_\Sigma \mathcal{L}) \boldsymbol{\alpha} \\ \nabla_M \mathcal{L} &= \nabla_\Sigma \mathcal{L}\end{aligned}$$

By plugging the gradients into the update in (14) and then re-expressing the update in terms of $\boldsymbol{\mu}$, $\boldsymbol{\alpha}$, $\boldsymbol{\Sigma}$, we obtain the natural gradient update in terms of $\boldsymbol{\mu}$, $\boldsymbol{\alpha}$, and $\boldsymbol{\Sigma}$.

$$\boldsymbol{\Sigma}^{-1} \leftarrow \boldsymbol{\Sigma}^{-1} - 2\beta \nabla_\Sigma \mathcal{L} \quad (58)$$

$$\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} + \beta \boldsymbol{\Sigma} (2\nabla_\mu \mathcal{L} - \nabla_\alpha \mathcal{L}) \quad (59)$$

$$\boldsymbol{\alpha} \leftarrow \boldsymbol{\alpha} + \beta \boldsymbol{\Sigma} (\nabla_\alpha \mathcal{L} - \nabla_\mu \mathcal{L}) \quad (60)$$

Recall that the lower bound is

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})} \left[-\sum_{n=1}^N f_n(\mathbf{z}) + \underbrace{\log \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})}_{\text{prior}} - \underbrace{\log q(\mathbf{z}|\boldsymbol{\lambda})}_{\text{entropy}} \right],$$

For the prior term, there is a closed-form expression for gradient computation.

$$\mathbb{E}_{q(\mathbf{z})} [\log \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})] = -\frac{d}{2} \log(2\pi) + \frac{d\delta}{2} - \frac{\delta}{2} (2\boldsymbol{\alpha}^T \boldsymbol{\alpha} + 2\boldsymbol{\mu}^T \boldsymbol{\alpha} + \text{Tr}(\boldsymbol{\Sigma}) + \boldsymbol{\mu}^T \boldsymbol{\mu})$$

It is easy to show that the gradients about the prior term are

$$\begin{aligned}\mathbf{g}_\mu^{\text{prior}} &= -\delta (\boldsymbol{\mu} + \boldsymbol{\alpha}) \\ \mathbf{g}_\alpha^{\text{prior}} &= -\delta (2\boldsymbol{\alpha} + \boldsymbol{\mu}) \\ \mathbf{g}_\Sigma^{\text{prior}} &= -\frac{\delta}{2} \mathbf{I}\end{aligned}$$

For the entropy term, it can be expressed as follows.

$$\mathbb{E}_{q(\mathbf{z})} [-\log q(\mathbf{z})] = \frac{1}{2} \left\{ \log \det(2\pi\boldsymbol{\Sigma}) + \log \left(\frac{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}{2\pi} \right) - \frac{1}{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}} + (d+1) \right\} - \mathbb{E}_{\text{Exp}(w|1)\mathcal{N}(z_2|w\sqrt{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}, 1)} \left[\log \phi \left(z_2 - \frac{1}{\sqrt{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}} \right) \right] \quad (61)$$

Using the expression in (61), the gradients of the entropy term are given as follows:

$$\mathbf{g}_\mu^{\text{entropy}} = 0 \quad (62)$$

$$\mathbf{g}_\alpha^{\text{entropy}} = \frac{\boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}} + \frac{\boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}{(\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha})^2} - \mathbb{E}_{\text{Exp}(w|1)\mathcal{N}(z_2|0,1)} \left[\frac{\exp \left(-\frac{t^2}{2} - \log \phi(t) \right)}{\sqrt{2\pi}} \left(w + \frac{1}{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}} \right) \left(\frac{\boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}{\sqrt{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}} \right) \right] \quad (63)$$

$$\mathbf{g}_\Sigma^{\text{entropy}} = \frac{\boldsymbol{\Sigma}^{-1}}{2} - \frac{\boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1}}{2\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}} - \frac{\boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1}}{2(\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha})^2} + \mathbb{E}_{\text{Exp}(w|1)\mathcal{N}(z_2|0,1)} \left[\frac{\exp \left(-\frac{t^2}{2} - \log \phi(t) \right)}{\sqrt{2\pi}} \left(w + \frac{1}{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}} \right) \left(\frac{\boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} \boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1}}{2\sqrt{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}} \right) \right] \quad (64)$$

where $t = z_2 + \frac{w\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} - 1}{\sqrt{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}}}$ and the expectations involve 2d integrations, which can be computed by Gauss-Hermite quadrature and Gauss-Laguerre quadrature.

The remaining thing is to compute the gradients about $\mathbb{E}_{q(z)} [f_n(\mathbf{z})]$. To compute the gradients, the reparametrization trick can be used. However, we can use the extended Bonnet's and Price's theorems (Lin et al., 2019). Assuming that $f_n(\mathbf{z})$ satisfies the assumptions needed for these two theorems, we obtain the following gradient expression:

$$\begin{aligned} \mathbf{g}_1^n &:= \nabla_{\mu} \mathbb{E}_{q(z)} [f_n(\mathbf{z})] = \mathbb{E}_{q(z)} [\nabla_z f_n(\mathbf{z})] \approx \nabla_z f_n(\mathbf{z}) \\ \mathbf{g}_2^n &:= \nabla_{\alpha} \mathbb{E}_{q(z)} [f_n(\mathbf{z})] \\ &= \mathbb{E}_{q(z)} [u(\mathbf{z}) \nabla_z f_n(\mathbf{z})] + v \mathbb{E}_{\mathcal{N}(\hat{\mathbf{z}}|\mu, \Sigma)} [\nabla_{\hat{\mathbf{z}}} f_n(\hat{\mathbf{z}})] \approx u(\mathbf{z}) \nabla_z f_n(\mathbf{z}) + v \nabla_{\hat{\mathbf{z}}} f_n(\hat{\mathbf{z}}) \\ &= \mathbb{E}_{q(w, z)} [w \nabla_z f_n(\mathbf{z})] \approx w \nabla_z f_n(\mathbf{z}) \\ \mathbf{g}_3^n &:= 2 \nabla_{\Sigma} \mathbb{E}_{q(z)} [f_n(\mathbf{z})] = \mathbb{E}_{q(z)} [\nabla_z^2 f_n(\mathbf{z})] \approx \nabla_z^2 f_n(\mathbf{z}) \end{aligned}$$

where $v = \frac{1}{(\alpha^T \Sigma^{-1} \alpha)}$, $u(\mathbf{z}) := \frac{(\mathbf{z} - \mu)^T \Sigma^{-1} \alpha - 1}{\alpha^T \Sigma^{-1} \alpha}$, and $w \sim \text{Exp}(w|1)$, $\hat{\mathbf{z}} \sim \mathcal{N}(\hat{\mathbf{z}}|\mu, \Sigma)$, $\mathbf{z} = \hat{\mathbf{z}} + w\alpha$.

Putting together, we can express the gradients in the following form, which will be used for deriving the extended variational Adam update.

$$\begin{aligned} \nabla_{\mu} \mathcal{L}(\lambda) &= - \sum_{n=1}^N \left(\underbrace{\mathbf{g}_1^n - \mathbf{g}_{\mu}^{\text{entropy}}/N}_0 \right) \underbrace{-\delta(\mu + \alpha)}_{\mathbf{g}_{\mu}^{\text{prior}}} \\ &= - \sum_{n=1}^N \underbrace{\mathbf{g}_1^n}_{:= \mathbf{g}_{\mu}^n} - \delta(\mu + \alpha) \end{aligned} \quad (65)$$

$$\nabla_{\alpha} \mathcal{L}(\lambda) = - \sum_{n=1}^N \left(\underbrace{\mathbf{g}_2^n - \mathbf{g}_{\alpha}^{\text{entropy}}/N}_{:= \mathbf{g}_{\alpha}^n} \right) \underbrace{-\delta(2\alpha + \mu)}_{\mathbf{g}_{\alpha}^{\text{prior}}} \quad (66)$$

$$\begin{aligned} \nabla_{\Sigma} \mathcal{L}(\lambda) &= -\frac{1}{2} \sum_{n=1}^N (\mathbf{g}_3^n - 2\mathbf{g}_{\Sigma}^{\text{entropy}}/N) \underbrace{-\frac{\delta}{2}\mathbf{I} + \frac{1}{2}(\Sigma^{-1} - \Sigma^{-1})}_{\mathbf{g}_{\Sigma}^{\text{prior}}} \\ &= -\frac{1}{2} \sum_{n=1}^N \left(\underbrace{\mathbf{g}_3^n - 2\mathbf{g}_{\Sigma}^{\text{entropy}}/N + \Sigma^{-1}/N}_{:= \mathbf{g}_{\Sigma}^n} \right) - \frac{\delta}{2}\mathbf{I} + \frac{1}{2}\Sigma^{-1} \end{aligned} \quad (67)$$

Similarly, for stochastic approximation, we can sub-sampling a data point n and use MC samples to approximate \mathbf{g}_1^n , \mathbf{g}_2^n , and \mathbf{g}_3^n . Plugging these stochastic gradients into (58)-(60), we obtain the NGD update:

$$\begin{aligned} \Sigma^{-1} &\leftarrow (1 - \beta)\Sigma^{-1} + \beta(\delta\mathbf{I} + N\mathbf{g}_{\Sigma}^n) \\ \mu &\leftarrow \mu - \beta\Sigma(N(2\mathbf{g}_{\mu}^n - \mathbf{g}_{\alpha}^n) + \delta\mu) \\ \alpha &\leftarrow \alpha - \beta\Sigma(N(\mathbf{g}_{\alpha}^n - \mathbf{g}_{\mu}^n) + \delta\alpha) \end{aligned}$$

G. Multivariate Normal Inverse-Gaussian Distribution

We consider the following mixture distribution (Barndorff-Nielsen, 1997), which is a Gaussian variance-mean mixture distribution. For simplicity, we assume λ is known.

$$q(w, \mathbf{z}) = \mathcal{N}(\mathbf{z}|\mu + w\alpha, w\Sigma) \text{InvGauss}(w|1, \lambda)$$

where $\text{InvGauss}(w|1, \lambda) = \left(\frac{\lambda}{2\pi w^3}\right)^{\frac{1}{2}} \exp\left\{-\frac{\lambda}{2}(w + w^{-1}) + \lambda\right\}$ denotes the inverse Gaussian distribution.

Lemma 15 *The marginal distribution is*

$$q(\mathbf{z}) = \frac{\lambda^{\frac{1}{2}}}{(2\pi)^{\frac{d+1}{2}}} \det(\Sigma)^{-1/2} \exp\left[(\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} \boldsymbol{\alpha} + \lambda\right] \frac{2\mathcal{K}_{\frac{d+1}{2}}\left(\sqrt{(\boldsymbol{\alpha}^T \Sigma^{-1} \boldsymbol{\alpha} + \lambda) \left((\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \lambda\right)}\right)}{\left(\sqrt{\frac{\boldsymbol{\alpha}^T \Sigma^{-1} \boldsymbol{\alpha} + \lambda}{(\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \lambda}}\right)^{\frac{-d-1}{2}}},$$

where $\mathcal{K}_v(x)$ denotes the modified Bessel function of the second kind.

Proof: By definition, we can compute the marginal distribution as follows.

$$\begin{aligned} q(\mathbf{z}) &= \int_0^{+\infty} q(\mathbf{z}|w)q(w)dw \\ &= \int_0^{+\infty} \det(2\pi w\Sigma)^{-1/2} \exp\left\{-\frac{1}{2}\left[(\mathbf{z} - \boldsymbol{\mu} - w\boldsymbol{\alpha})^T (w\Sigma)^{-1} (\mathbf{z} - \boldsymbol{\mu} - w\boldsymbol{\alpha})\right]\right\} \left(\frac{\lambda}{2\pi w^3}\right)^{1/2} \exp\left[-\frac{\lambda}{2}\left(w + \frac{1}{w}\right) + \lambda\right] dw \\ &\quad \text{By grouping all terms related to } w \text{ together, we have} \\ &= \frac{\lambda^{\frac{1}{2}}}{(2\pi)^{\frac{d+1}{2}}} \det(\Sigma)^{-1/2} \exp\left[(\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} \boldsymbol{\alpha} + \lambda\right] \int_0^{+\infty} w^{-\frac{d+3}{2}} \exp\left\{-\frac{1}{2}\left[w(\boldsymbol{\alpha}^T \Sigma^{-1} \boldsymbol{\alpha} + \lambda) + \frac{(\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \lambda}{w}\right]\right\} dw \end{aligned}$$

By completing a generalized inverse Gaussian form, we have

$$= \frac{\lambda^{\frac{1}{2}}}{(2\pi)^{\frac{d+1}{2}}} \det(\Sigma)^{-1/2} \exp\left[(\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} \boldsymbol{\alpha} + \lambda\right] \frac{2\mathcal{K}_{\frac{-d-1}{2}}\left(\sqrt{(\boldsymbol{\alpha}^T \Sigma^{-1} \boldsymbol{\alpha} + \lambda) \left((\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \lambda\right)}\right)}{\left(\sqrt{\frac{\boldsymbol{\alpha}^T \Sigma^{-1} \boldsymbol{\alpha} + \lambda}{(\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \lambda}}\right)^{\frac{-d-1}{2}}}$$

We obtain the last step by the fact that $\mathcal{K}_v(x) = \mathcal{K}_{-v}(x)$

$$= \frac{\lambda^{\frac{1}{2}}}{(2\pi)^{\frac{d+1}{2}}} \det(\Sigma)^{-1/2} \exp\left[(\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} \boldsymbol{\alpha} + \lambda\right] \frac{2\mathcal{K}_{\frac{d+1}{2}}\left(\sqrt{(\boldsymbol{\alpha}^T \Sigma^{-1} \boldsymbol{\alpha} + \lambda) \left((\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \lambda\right)}\right)}{\left(\sqrt{\frac{\boldsymbol{\alpha}^T \Sigma^{-1} \boldsymbol{\alpha} + \lambda}{(\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \lambda}}\right)^{\frac{-d-1}{2}}}$$

□

Similarly, the sufficient statistics $\phi_z(\mathbf{z}, w) = \{\mathbf{z}/w, \mathbf{z}, \mathbf{z}\mathbf{z}^T/w\}$ and the natural parameter $\boldsymbol{\lambda}_z = \{\Sigma^{-1}\boldsymbol{\mu}, \Sigma^{-1}\boldsymbol{\alpha}, -\frac{1}{2}\Sigma^{-1}\}$ can be read from $q(\mathbf{z}|w) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu} + w\boldsymbol{\alpha}, w\Sigma)$. The joint distribution $q(\mathbf{z}, w)$ is a conditional EF because $q(w)$ is an inverse Gaussian distribution with known parameters, which is a EF distribution. Likewise, we can show that the joint distribution is also a minimal conditional EF.

We can derive the expectation parameter shown below:

$$\begin{aligned} \mathbf{m} &:= \mathbb{E}_{\text{InvGauss}(w|1, \lambda), \mathcal{N}(z|\mu + w\alpha, w\Sigma)} [w^{-1}\mathbf{z}] = (1 + \lambda^{-1})\boldsymbol{\mu} + \boldsymbol{\alpha}, \\ \mathbf{m}_\alpha &:= \mathbb{E}_{\text{InvGauss}(w|1, \lambda), \mathcal{N}(z|\mu + w\alpha, w\Sigma)} [\mathbf{z}] = \boldsymbol{\mu} + \boldsymbol{\alpha}, \\ \mathbf{M} &:= \mathbb{E}_{\text{InvGauss}(w|1, \lambda), \mathcal{N}(z|\mu + w\alpha, w\Sigma)} [w^{-1}\mathbf{z}\mathbf{z}^T] = (1 + \lambda^{-1})\boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\alpha}\boldsymbol{\alpha}^T + \boldsymbol{\mu}\boldsymbol{\alpha}^T + \boldsymbol{\alpha}\boldsymbol{\mu}^T + \Sigma \end{aligned}$$

The sufficient statistics, natural parameters, and expectation parameters are summarized below:

$$\begin{bmatrix} \mathbf{z}/w \\ \mathbf{z} \\ \mathbf{z}\mathbf{z}^T/w \end{bmatrix} \quad \begin{bmatrix} \Sigma^{-1}\boldsymbol{\mu} \\ \Sigma^{-1}\boldsymbol{\alpha} \\ -\frac{1}{2}\Sigma^{-1} \end{bmatrix} \quad \begin{bmatrix} (1 + \lambda^{-1})\boldsymbol{\mu} + \boldsymbol{\alpha} \\ \boldsymbol{\mu} + \boldsymbol{\alpha} \\ (1 + \lambda^{-1})\boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\alpha}\boldsymbol{\alpha}^T + \boldsymbol{\mu}\boldsymbol{\alpha}^T + \boldsymbol{\alpha}\boldsymbol{\mu}^T + \Sigma \end{bmatrix}$$

G.1. Derivation of the NGD Update

We consider the variational approximation using the normal inverse Gaussian distribution $q(\mathbf{z}|\boldsymbol{\lambda})$. We consider the same model as Appendix B.1 with a Gaussian-prior $\mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})$ on \mathbf{z} . The lower bound is defined as below:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\lambda}) &= \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})} \left[\sum_{n=1}^N \log p(\mathcal{D}_n|\mathbf{z}) + \log \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I}) - \log q(\mathbf{z}|\boldsymbol{\lambda}) \right] \\ &= \mathbb{E}_{q(\mathbf{z})} [-h(\mathbf{z})], \text{ where } h(\mathbf{z}) := - \left[\log \frac{\mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})}{q(\mathbf{z})} + \sum_n \log p(\mathcal{D}_n|\mathbf{z}) \right].\end{aligned}$$

Our goal is to compute the gradient of this ELBO with respect to the expectation parameters.

G.2. Natural Gradient for $q(\mathbf{z}|w)$

Likewise, we do not need to compute these gradients with respect to the expectation parameters directly. The gradients can be computed in terms of $\{\boldsymbol{\mu}, \boldsymbol{\alpha}, \boldsymbol{\Sigma}\}$. Similarly, by the chain rule, we can express the following gradients with respect to the expectation parameters in terms of the gradients with respect to $\boldsymbol{\mu}$, $\boldsymbol{\alpha}$, and $\boldsymbol{\Sigma}$.

$$\begin{aligned}\nabla_m \mathcal{L} &= \lambda \nabla_{\boldsymbol{\mu}} \mathcal{L} - \lambda \nabla_{\boldsymbol{\alpha}} \mathcal{L} - 2 (\nabla_{\boldsymbol{\Sigma}} \mathcal{L}) \boldsymbol{\mu} \\ \nabla_{m_{\alpha}} \mathcal{L} &= (1 + \lambda) \nabla_{\boldsymbol{\alpha}} \mathcal{L} - \lambda \nabla_{\boldsymbol{\mu}} \mathcal{L} - 2 (\nabla_{\boldsymbol{\Sigma}} \mathcal{L}) \boldsymbol{\alpha} \\ \nabla_M \mathcal{L} &= \nabla_{\boldsymbol{\Sigma}} \mathcal{L}\end{aligned}$$

By plugging the gradients into the update in (14) and then re-expressing the update in terms of $\boldsymbol{\mu}$, $\boldsymbol{\alpha}$, $\boldsymbol{\Sigma}$, we obtain the natural gradient update in terms of $\boldsymbol{\mu}$, $\boldsymbol{\alpha}$, and $\boldsymbol{\Sigma}$.

$$\boldsymbol{\Sigma}^{-1} \leftarrow \boldsymbol{\Sigma}^{-1} - 2\beta \nabla_{\boldsymbol{\Sigma}} \mathcal{L} \quad (68)$$

$$\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} + \beta \boldsymbol{\Sigma} (\lambda \nabla_{\boldsymbol{\mu}} \mathcal{L} - \lambda \nabla_{\boldsymbol{\alpha}} \mathcal{L}) \quad (69)$$

$$\boldsymbol{\alpha} \leftarrow \boldsymbol{\alpha} + \beta \boldsymbol{\Sigma} ((1 + \lambda) \nabla_{\boldsymbol{\alpha}} \mathcal{L} - \lambda \nabla_{\boldsymbol{\mu}} \mathcal{L}) \quad (70)$$

We can compute gradients with respect to $\boldsymbol{\mu}$, $\boldsymbol{\alpha}$, and $\boldsymbol{\Sigma}$ by the extended Bonnet's and Price's theorem (Lin et al., 2019). In Lin et al. (2019), they discuss the conditions of the target function $h(\mathbf{z})$ when it comes to applying these theorems.

$$\begin{aligned}\nabla_{\boldsymbol{\mu}} \mathcal{L}(\boldsymbol{\lambda}) &= - \mathbb{E}_{q(\mathbf{z})} [\nabla_{\mathbf{z}} h(\mathbf{z})] \approx - \nabla_{\mathbf{z}} h(\mathbf{z}) \\ \nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\lambda}) &= - \mathbb{E}_{q(w, \mathbf{z})} [w \nabla_{\mathbf{z}} h_n(\mathbf{z})] \approx - w \nabla_{\mathbf{z}} h_n(\mathbf{z}) \\ &= - \mathbb{E}_{q(\mathbf{z})} [u(\mathbf{z}) \nabla_{\mathbf{z}} h_n(\mathbf{z})] \approx - u(\mathbf{z}) \nabla_{\mathbf{z}} h_n(\mathbf{z}) \\ \nabla_{\boldsymbol{\Sigma}} \mathcal{L}(\boldsymbol{\lambda}) &= - \frac{1}{2} \mathbb{E}_{q(\mathbf{z})} [w \nabla_{\mathbf{z}}^2 h(\mathbf{z})] \approx - \frac{w}{2} \nabla_{\mathbf{z}}^2 h(\mathbf{z}). \\ &= - \frac{1}{2} \mathbb{E}_{q(\mathbf{z})} [u(\mathbf{z}) \nabla_{\mathbf{z}}^2 h(\mathbf{z})] \approx - \frac{u(\mathbf{z})}{2} \nabla_{\mathbf{z}}^2 h(\mathbf{z}).\end{aligned}$$

where $u(\mathbf{z}) := \sqrt{\frac{(\mathbf{z}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z}-\boldsymbol{\mu}) + \lambda}{\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} + \lambda}} \frac{\mathcal{K}_{\frac{d-1}{2}} \left(\sqrt{(\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} + \lambda) ((\mathbf{z}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z}-\boldsymbol{\mu}) + \lambda)} \right)}{\mathcal{K}_{\frac{d+1}{2}} \left(\sqrt{(\boldsymbol{\alpha}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} + \lambda) ((\mathbf{z}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z}-\boldsymbol{\mu}) + \lambda)} \right)}$ and $w \sim \text{InvGauss}(w|1, \lambda)$, $\mathbf{z} \sim \mathcal{N}(\mathbf{z}|\boldsymbol{\mu} + w\boldsymbol{\alpha}, w\boldsymbol{\Sigma})$.

Directly calculating the ratio between the modified Bessel functions of the second kind is expensive and numerically unstable when v is large. However, the ratio between consecutive order has a tight and algebraic bound (Ruiz-Antolín & Segura, 2016) as given below. When $v \in \mathcal{R}$ and $v \geq \frac{1}{2}$, the bound of the ratio is

$$D_{2v-1}(v, x) \leq \frac{\mathcal{K}_{v-1}(x)}{\mathcal{K}_v(x)} \leq D_0(v, x)$$

where function $D_\alpha(v, x)$ is defined as below.

$$D_\alpha(v, x) := \frac{x}{\psi_\alpha(v, x) + \sqrt{(\psi_\alpha(v, x))^2 + x^2}}$$

$$\psi_\alpha(v, x) := (v - \frac{1}{2}) - \frac{\tau_\alpha(v)}{2\sqrt{(\tau_\alpha(v))^2 + x^2}}, \quad \tau_\alpha(v) := v - \frac{\alpha + 1}{2}$$

For natural number $v \in \mathbb{N}$, a tighter bound (see Eq (3.10) at Yang & Chu (2017)) with higher computation cost can be used, where we make use of the following relationship due to Eq (72) and (73).

$$\frac{\mathcal{K}_{v-1}(x)}{\mathcal{K}_v(x)} = \frac{\mathcal{K}_{v+1}(x)}{\mathcal{K}_v(x)} - \frac{2v}{x}$$

To compute the ratio, we propose to use the following approximation when $v \geq \frac{1}{2}$. A similar approach to approximate the ratio between two modified Bessel functions of the first kind is used in Oh et al. (2019) and Kumar & Tsvetkov (2018).

$$\frac{\mathcal{K}_{v-1}(x)}{\mathcal{K}_v(x)} \approx \frac{D_{2v-1}(v, x) + D_0(v, x)}{2} \quad (71)$$

Now, we discuss how to compute $\nabla_x \log \mathcal{K}_v(x)$ and $\nabla_x^2 \log \mathcal{K}_v(x)$. The first term appears when we compute $\nabla_z h(\mathbf{z})$. Similarly, the second term appears when we compute $\nabla_z^2 h(\mathbf{z})$.

First, we make use of the recurrence forms of the modified Bessel function of the second kind for $v \in \mathcal{R}$ (see page 20 at Culham (2004)).

$$\nabla_x \mathcal{K}_v(x) = -\mathcal{K}_{v-1}(x) - \frac{v}{x} \mathcal{K}_v(x) \quad (72)$$

$$\nabla_x \mathcal{K}_v(x) = \frac{v}{x} \mathcal{K}_v(x) - \mathcal{K}_{v+1}(x) \quad (73)$$

Using these recurrence forms, we have

$$\frac{\nabla_x \mathcal{K}_v(x)}{\mathcal{K}_v(x)} = -\frac{\mathcal{K}_{v-1}(x)}{\mathcal{K}_v(x)} - \frac{v}{x} \quad (74)$$

Furthermore, we have the following result due to the the recurrence forms.

$$\begin{aligned} \nabla_x^2 \mathcal{K}_v(x) &= \nabla_x \left[\underbrace{-\mathcal{K}_{v-1}(x) - \frac{v}{x} \mathcal{K}_v(x)}_{\nabla_x \mathcal{K}_v(x)} \right] \\ &= -\nabla_x \mathcal{K}_{v-1}(x) - \frac{v}{x} \nabla_x \mathcal{K}_v(x) + \frac{v}{x^2} \mathcal{K}_v(x) \\ &= -\left(\underbrace{\frac{v-1}{x} \mathcal{K}_{v-1}(x) - \mathcal{K}_v(x)}_{\nabla_x \mathcal{K}_{v-1}(x)} \right) - \frac{v}{x} \nabla_x \mathcal{K}_v(x) + \frac{v}{x^2} \mathcal{K}_v(x) \\ &= -\frac{v-1}{x} \mathcal{K}_{v-1}(x) + \frac{v+x^2}{x^2} \mathcal{K}_v(x) - \frac{v}{x} \nabla_x \mathcal{K}_v(x) \\ &= -\frac{v-1}{x} \mathcal{K}_{v-1}(x) + \frac{v+x^2}{x^2} \mathcal{K}_v(x) - \frac{v}{x} \left(\underbrace{-\mathcal{K}_{v-1}(x) - \frac{v}{x} \mathcal{K}_v(x)}_{\nabla_x \mathcal{K}_v(x)} \right) \\ &= \frac{1}{x} \mathcal{K}_{v-1}(x) + \frac{v+x^2+v^2}{x^2} \mathcal{K}_v(x) \end{aligned}$$

which implies that

$$\frac{\nabla_x^2 \mathcal{K}_v(x)}{\mathcal{K}_v(x)} = \frac{1}{x} \frac{\mathcal{K}_{v-1}(x)}{\mathcal{K}_v(x)} + \frac{v + x^2 + v^2}{x^2} \quad (75)$$

Using Eq (74) and (75), we have

$$\begin{aligned} \nabla_x \log \mathcal{K}_v(x) &= \frac{\nabla_x \mathcal{K}_v(x)}{\mathcal{K}_v(x)} \\ &= -\frac{\mathcal{K}_{v-1}(x)}{\mathcal{K}_v(x)} - \frac{v}{x} \\ \nabla_x^2 \log \mathcal{K}_v(x) &= \nabla_x \left[\frac{\nabla_x \mathcal{K}_v(x)}{\mathcal{K}_v(x)} \right] \\ &= \frac{\nabla_x^2 \mathcal{K}_v(x)}{\mathcal{K}_v(x)} - \left(\frac{\nabla_x \mathcal{K}_v(x)}{\mathcal{K}_v(x)} \right)^2 \\ &= \frac{1}{x} \frac{\mathcal{K}_{v-1}(x)}{\mathcal{K}_v(x)} + \frac{v + x^2 + v^2}{x^2} - \left(-\frac{\mathcal{K}_{v-1}(x)}{\mathcal{K}_v(x)} - \frac{v}{x} \right)^2 \\ &= \frac{1 - 2v}{x} \left(\frac{\mathcal{K}_{v-1}(x)}{\mathcal{K}_v(x)} \right) - \left(\frac{\mathcal{K}_{v-1}(x)}{\mathcal{K}_v(x)} \right)^2 + \frac{v}{x^2} + 1 \end{aligned}$$

where the ratio can be approximated by Eq (71).

H. Multivariate Symmetric Normal Inverse-Gaussian Distribution

The symmetric normal inverse-Gaussian distribution is a scale mixture distribution. A difference between the distribution at Appendix G is shown in red. Such difference allows this distribution to have heavy tails.

$$q(w, \mathbf{z}) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, w^{-1} \boldsymbol{\Sigma}) \text{InvGauss}(w | 1, \lambda)$$

where we assume λ is fixed for simplicity and $\text{InvGauss}(w | 1, \lambda) = \left(\frac{\lambda}{2\pi w^3} \right)^{\frac{1}{2}} \exp \left\{ -\frac{\lambda}{2} (w + w^{-1}) + \lambda \right\}$.

Similarly, we can show that the marginal distribution is

$$q(\mathbf{z}) = \frac{\lambda^{\frac{1}{2}}}{(2\pi)^{\frac{d+1}{2}} \det(\boldsymbol{\Sigma})^{-1/2} \exp(\lambda)} \frac{2\mathcal{K}_{\frac{d-1}{2}} \left(\sqrt{\lambda \left((\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \lambda \right)} \right)}{\left(\sqrt{\frac{(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \lambda}{\lambda}} \right)^{\frac{d-1}{2}}}$$

Furthermore, the mixture distribution is a minimal conditional EF distribution. The sufficient statistics, natural parameters, and expectation parameters are summarized below:

$$\begin{bmatrix} w\mathbf{z} \\ w\mathbf{z}\mathbf{z}^T \end{bmatrix} \quad \begin{bmatrix} \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ -\frac{1}{2}\boldsymbol{\Sigma}^{-1} \end{bmatrix} \quad \begin{bmatrix} \boldsymbol{\mu} \\ \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma} \end{bmatrix}$$

Likewise, we can derive the natural-gradient update in terms of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ as shown below.

$$\begin{aligned} \boldsymbol{\Sigma}^{-1} &\leftarrow \boldsymbol{\Sigma}^{-1} - 2\beta \nabla_{\boldsymbol{\Sigma}} \mathcal{L} \\ \boldsymbol{\mu} &\leftarrow \boldsymbol{\mu} + \beta \boldsymbol{\Sigma} \nabla_{\boldsymbol{\mu}} \mathcal{L} \end{aligned}$$

Now, we discuss the gradient computation of the following lower bound.

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})} \left[\sum_{n=1}^N \underbrace{\log p(\mathcal{D}_n | \mathbf{z})}_{:-f_n(\mathbf{z})} + \log \mathcal{N}(\mathbf{z} | \mathbf{0}, \delta^{-1} \mathbf{I}) - \log q(\mathbf{z} | \boldsymbol{\lambda}) \right]$$

For the prior term, we have

$$\mathbb{E}_{q(z)} [\log \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})] = -\frac{d}{2} \log(2\pi) + \frac{d\delta}{2} - \frac{\delta}{2} (\boldsymbol{\mu}^T \boldsymbol{\mu} + (1 + \lambda^{-1}) \text{Tr}(\boldsymbol{\Sigma})) \quad (76)$$

Due to Eq (76), the closed-form gradients are given below.

$$\begin{aligned} \mathbf{g}_{\boldsymbol{\mu}}^{\text{prior}} &= -\delta \boldsymbol{\mu} \\ \mathbf{g}_{\boldsymbol{\Sigma}}^{\text{prior}} &= -\frac{\delta(1 + \lambda^{-1})}{2} \mathbf{I} \end{aligned}$$

In this example, we can re-express the entropy term as below.

$$\mathbb{E}_{q(z)} [-\log q(\mathbf{z})] = \frac{(d+1) \log(2\pi)}{2} + \frac{\log \det(\boldsymbol{\Sigma})}{2} - \lambda - \frac{d+1}{4} \log(\lambda) - \mathbb{E}_{\text{InvGauss}(w|1, \lambda)} \chi^2(z_1|d) \left[\log \left(\frac{2\mathcal{K}_{\frac{d-1}{2}} \left(\sqrt{\lambda} (w^{-1}z_1 + \lambda) \right)}{(w^{-1}z_1 + \lambda)^{\frac{d-1}{4}}} \right) \right] \quad (77)$$

where $\chi^2(z_1|d)$ denotes the chi-squared distribution with d degrees of freedom. The expectation can be computed by the inverse Gaussian quadrature (Choi et al., 2018) and the generalized Gauss-Laguerre quadrature.

By Eq (77), the closed-form gradients of the entropy term are shown below.

$$\begin{aligned} \mathbf{g}_{\boldsymbol{\mu}}^{\text{entropy}} &= 0 \\ \mathbf{g}_{\boldsymbol{\Sigma}}^{\text{entropy}} &= \frac{\boldsymbol{\Sigma}^{-1}}{2} \end{aligned}$$

The remaining step is to compute the gradients about $\mathbb{E}_{q(z)} [f_n(\mathbf{z})]$. To compute the gradients, we use the extended the Bonnet's and Price's theorems (Lin et al., 2019). Assuming that $f_n(\mathbf{z})$ satisfies the assumptions needed for these two theorems, we obtain the following gradient expression:

$$\begin{aligned} \mathbf{g}_1^n &:= \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(z)} [f_n(\mathbf{z})] = \mathbb{E}_{q(z)} [\nabla_z f_n(\mathbf{z})] \approx \nabla_z f_n(\mathbf{z}) \\ \mathbf{g}_2^n &:= 2\nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{q(z)} [f_n(\mathbf{z})] \\ &= \mathbb{E}_{q(z)} [u(\mathbf{z}) \nabla_z^2 f_n(\mathbf{z})] \approx u(\mathbf{z}) \nabla_z^2 f_n(\mathbf{z}) \\ &= \mathbb{E}_{q(w, z)} [w \nabla_z^2 f_n(\mathbf{z})] \approx w^{-1} \nabla_z^2 f_n(\mathbf{z}) \end{aligned}$$

where $u(\mathbf{z}) := \sqrt{\frac{(\mathbf{z}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z}-\boldsymbol{\mu}) + \lambda}{\lambda}} \frac{\mathcal{K}_{\frac{d-3}{2}} \left(\sqrt{\lambda} ((\mathbf{z}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z}-\boldsymbol{\mu}) + \lambda) \right)}{\mathcal{K}_{\frac{d-1}{2}} \left(\sqrt{\lambda} ((\mathbf{z}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z}-\boldsymbol{\mu}) + \lambda) \right)}$, where the ratio between the Bessel functions

can be approximated by Eq (71) when $d \geq 2$, and $w \sim \text{InvGauss}(w|1, \lambda)$, $\mathbf{z} \sim \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, w^{-1}\boldsymbol{\Sigma})$.

Putting together, we can express the gradients in the following form:

$$\nabla_{\boldsymbol{\mu}} \mathcal{L}(\boldsymbol{\lambda}) = -\sum_{n=1}^N \underbrace{\mathbf{g}_1^n}_{:=\mathbf{g}_{\boldsymbol{\mu}}^n} - \delta \boldsymbol{\mu} \quad (78)$$

$$\nabla_{\boldsymbol{\Sigma}} \mathcal{L}(\boldsymbol{\lambda}) = -\frac{1}{2} \sum_{n=1}^N \underbrace{\mathbf{g}_2^n}_{:=\mathbf{g}_{\boldsymbol{\Sigma}}^n} - \frac{\delta(1 + \lambda^{-1})}{2} \mathbf{I} + \frac{1}{2} \boldsymbol{\Sigma}^{-1} \quad (79)$$

Similarly, for stochastic approximation, we can sub-sampling a data point n and use MC samples to approximate \mathbf{g}_1^n and \mathbf{g}_2^n . Plugging these stochastic gradients into

$$\begin{aligned} \boldsymbol{\Sigma}^{-1} &\leftarrow (1 - \beta) \boldsymbol{\Sigma}^{-1} + \beta (\delta (1 + \lambda^{-1}) \mathbf{I} + N \mathbf{g}_{\boldsymbol{\Sigma}}^n) \\ \boldsymbol{\mu} &\leftarrow \boldsymbol{\mu} - \beta \boldsymbol{\Sigma} (N \mathbf{g}_{\boldsymbol{\mu}}^n + \delta \boldsymbol{\mu}) \end{aligned}$$

I. Matrix-Variate Gaussian Distribution

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

Lemma 16 *Matrix Gaussian distribution is a member of the multi-linear exponential family.*

Proof: Let $\Lambda_1 = \mathbf{W}$, $\Lambda_2 = \mathbf{U}^{-1}$, and $\Lambda_3 = \mathbf{V}^{-1}$. The distribution on $\mathbf{Z} \in \mathcal{R}^{d \times p}$ can be expressed as follows.

$$\begin{aligned} \mathcal{MN}(\mathbf{Z}|\mathbf{W}, \mathbf{U}, \mathbf{V}) &= (2\pi)^{-dp/2} \exp \left[-\frac{1}{2} \text{Tr} \left(\mathbf{V}^{-1} (\mathbf{Z} - \mathbf{W})^T \mathbf{U}^{-1} (\mathbf{Z} - \mathbf{W}) \right) - (d/2 \log \text{Det}(\mathbf{V}) + p/2 \log \text{Det}(\mathbf{U})) \right] \\ &= (2\pi)^{-dp/2} \exp \left\{ \text{Tr} \left(\Lambda_3 \left(-\frac{1}{2} \mathbf{Z} + \Lambda_1 \right)^T \Lambda_2 \mathbf{Z} \right) \right. \\ &\quad \left. - \frac{1}{2} \left[\text{Tr} \left(\Lambda_3 \Lambda_1^T \Lambda_2 \Lambda_1 \right) + d \log \text{Det}(\Lambda_3) + p \log \text{Det}(\Lambda_2) \right] \right\}. \end{aligned}$$

The function $\text{Tr} \left(\Lambda_3 \left(-\frac{1}{2} \mathbf{Z} + \Lambda_1 \right)^T \Lambda_2 \mathbf{Z} \right)$ is linear with respect each Λ_j given others. \square

We now derive the NGVI update using our new expectation parameterization. We can obtain function ϕ_1 , ϕ_2 , and ϕ_3 from the multi-linear function

$$f(\mathbf{Z}, \Lambda) := \text{Tr} \left(\Lambda_3 \left(-\frac{1}{2} \mathbf{Z} + \Lambda_1 \right)^T \Lambda_2 \mathbf{Z} \right).$$

For example, we can obtain function ϕ_1 from $f(\mathbf{Z}, \Lambda)$ as shown below:

$$f(\mathbf{Z}, \Lambda) = \langle \Lambda_1, \underbrace{\Lambda_2 \mathbf{Z} \Lambda_3}_{\phi_1(\mathbf{Z}, \Lambda_{-1})} \rangle - \underbrace{\frac{1}{2} \text{Tr} \left(\Lambda_3 \mathbf{Z}^T \Lambda_2 \mathbf{Z} \right)}_{r_1(\mathbf{Z}, \Lambda_{-1})}.$$

Similarly, we can obtain functions ϕ_2 and ϕ_3 . The corresponding expectation parameters of the Matrix Gaussian distribution can then be derived as below:

$$\begin{aligned} \mathbf{M}_1 &= \mathbb{E}_{\mathcal{MN}(\mathbf{Z}|\mathbf{W}, \mathbf{U}, \mathbf{V})} [\Lambda_2 \mathbf{Z} \Lambda_3] = \Lambda_2 \Lambda_1 \Lambda_3 \\ \mathbf{M}_2 &= \mathbb{E}_{\mathcal{MN}(\mathbf{Z}|\mathbf{W}, \mathbf{U}, \mathbf{V})} \left[-\frac{1}{2} \mathbf{Z} \Lambda_3 \mathbf{Z}^T + \mathbf{Z} \Lambda_3 \Lambda_1^T \right] = \frac{1}{2} \left(\Lambda_1 \Lambda_3 \Lambda_1^T - p \Lambda_2^{-1} \right) \\ \mathbf{M}_3 &= \mathbb{E}_{\mathcal{MN}(\mathbf{Z}|\mathbf{W}, \mathbf{U}, \mathbf{V})} \left[-\frac{1}{2} \mathbf{Z}^T \Lambda_2 \mathbf{Z} + \Lambda_1^T \Lambda_2 \mathbf{Z} \right] = \frac{1}{2} \left(\Lambda_1^T \Lambda_2 \Lambda_1 - d \Lambda_3^{-1} \right) \end{aligned}$$

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

$$\begin{aligned} \nabla_{M_1} \mathbb{E}_{q(\mathbf{Z}|\lambda)} [h(\mathbf{Z})] &= (\Lambda_2)^{-1} \nabla_{\mathbf{W}} \mathbb{E}_{\mathcal{MN}(\mathbf{Z}|\mathbf{W}, \mathbf{U}, \mathbf{V})} [h(\mathbf{Z})] (\Lambda_3)^{-1} \\ \nabla_{M_2} \mathbb{E}_{q(\mathbf{Z}|\lambda)} [h(\mathbf{Z})] &= \frac{-2}{p} \nabla_{\mathbf{U}} \mathbb{E}_{\mathcal{MN}(\mathbf{Z}|\mathbf{W}, \mathbf{U}, \mathbf{V})} [h(\mathbf{Z})] \\ \nabla_{M_3} \mathbb{E}_{q(\mathbf{Z}|\lambda)} [h(\mathbf{Z})] &= \frac{-2}{d} \nabla_{\mathbf{V}} \mathbb{E}_{\mathcal{MN}(\mathbf{Z}|\mathbf{W}, \mathbf{U}, \mathbf{V})} [h(\mathbf{Z})] \end{aligned}$$

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 $\mathbb{E}_{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})} [-\hat{h}(\mathbf{z})]$, where $\hat{h}(\mathbf{z}) = h(\mathbf{Z})$. We make use of the Bonnet's and Price's theorems (Oppen & Archambeau, 2009):

$$\begin{aligned} \nabla_{\boldsymbol{\mu}} \mathbb{E}_{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\hat{h}(\mathbf{z})] &= \mathbb{E}_{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\nabla_{\mathbf{z}} \hat{h}(\mathbf{z})] \\ \nabla_{\boldsymbol{\Sigma}} \mathbb{E}_{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\hat{h}(\mathbf{z})] &= \frac{1}{2} \mathbb{E}_{\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\nabla_{\mathbf{z}}^2 \hat{h}(\mathbf{z})] \end{aligned}$$

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} :

$$\begin{aligned}
 \nabla_W \mathbb{E}_{\mathcal{MN}(Z|W,U,V)} [h(\mathbf{Z})] &= \text{Mat} \left(\mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[\nabla_z \hat{h}(\mathbf{z}) \right] \right) \\
 &= \mathbb{E}_{\mathcal{MN}(Z|W,U,V)} [\nabla_Z h(\mathbf{Z})] \\
 \nabla_U \mathbb{E}_{\mathcal{MN}(Z|W,U,V)} [h(\mathbf{Z})] &= (\nabla_U \Sigma) \nabla_\Sigma \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[\nabla_z \hat{h}(\mathbf{z}) \right] \\
 &= \frac{1}{2} (\nabla_U \Sigma) \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[\nabla_z^2 \hat{h}(\mathbf{z}) \right] \\
 &\approx \frac{1}{2} (\nabla_U \Sigma) \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[\nabla_z \hat{h}(\mathbf{z}) \nabla_z \hat{h}(\mathbf{z})^T \right] \\
 &= \frac{1}{2} \mathbb{E}_{\mathcal{MN}(Z|W,U,V)} [\nabla_Z h(\mathbf{Z}) \mathbf{V} \nabla_Z h(\mathbf{Z})^T]
 \end{aligned} \tag{80}$$

$$\begin{aligned}
 \nabla_V \mathbb{E}_{\mathcal{MN}(Z|W,U,V)} [h(\mathbf{Z})] &= (\nabla_V \Sigma) \nabla_\Sigma \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[\nabla_z \hat{h}(\mathbf{z}) \right] \\
 &= \frac{1}{2} (\nabla_V \Sigma) \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[\nabla_z^2 \hat{h}(\mathbf{z}) \right] \\
 &\approx \frac{1}{2} (\nabla_V \Sigma) \mathbb{E}_{\mathcal{N}(z|\mu,\Sigma)} \left[\nabla_z \hat{h}(\mathbf{z}) \nabla_z \hat{h}(\mathbf{z})^T \right] \\
 &= \frac{1}{2} \mathbb{E}_{\mathcal{MN}(Z|W,U,V)} [\nabla_Z h(\mathbf{Z})^T \mathbf{U} \nabla_Z h(\mathbf{Z})] .
 \end{aligned} \tag{81}$$

To avoid computation of the Hessian, we have used the Gauss-Newton approximation (Khan et al., 2018) in Eq (80) and Eq. (81).

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

$$\begin{aligned}
 \mathbf{\Lambda}_1 &\leftarrow \mathbf{\Lambda}_1 - \beta_1 (\mathbf{\Lambda}_2)^{-1} \mathbb{E}_{\mathcal{MN}(Z|W,U,V)} [\nabla_Z h(\mathbf{Z})] (\mathbf{\Lambda}_3)^{-1} \\
 \mathbf{\Lambda}_2 &\leftarrow \mathbf{\Lambda}_2 + \beta_2 \mathbb{E}_{\mathcal{MN}(Z|W,U,V)} [\nabla_Z h(\mathbf{Z}) \mathbf{V} \nabla_Z h(\mathbf{Z})^T] \\
 \mathbf{\Lambda}_3 &\leftarrow \mathbf{\Lambda}_3 + \beta_2 \mathbb{E}_{\mathcal{MN}(Z|W,U,V)} [\nabla_Z h(\mathbf{Z})^T \mathbf{U} \nabla_Z h(\mathbf{Z})]
 \end{aligned}$$

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

$$\begin{aligned}
 \mathbf{W} &\leftarrow \mathbf{W} - \beta_1 \mathbf{U} \mathbb{E}_{\mathcal{MN}(Z|W,U,V)} [\nabla_Z h(\mathbf{Z})] \mathbf{V} \\
 (\mathbf{U})^{-1} &\leftarrow (\mathbf{U})^{-1} + \beta_2 \mathbb{E}_{\mathcal{MN}(Z|W,U,V)} [\nabla_Z h(\mathbf{Z}) \mathbf{V} \nabla_Z h(\mathbf{Z})^T] \\
 (\mathbf{V})^{-1} &\leftarrow (\mathbf{V})^{-1} + \beta_2 \mathbb{E}_{\mathcal{MN}(Z|W,U,V)} [\nabla_Z h(\mathbf{Z})^T \mathbf{U} \nabla_Z h(\mathbf{Z})]
 \end{aligned}$$

J. Extensions to Variational Adam

For simplicity, we consider a case when variational parameters of $q(w|\lambda_w)$ are fixed. Since λ_w is fixed, using the same derivation as Khan et al. (2018), we obtain the following natural-gradient update with the natural momentum ($0 \leq m < 1$).

$$\lambda_z^{t+1} = \frac{1}{1-m} \lambda_z^t - \frac{m}{1-m} \lambda_z^{t-1} + \frac{\beta}{1-m} \nabla_{m_z} \mathcal{L}(\lambda_z) \Big|_{\lambda_z = \lambda_z^t} \tag{82}$$

We assume the model prior is a Gaussian prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})$ to derive extensions of the variational Adam update, where the variational distribution is a Gaussian mixture distribution such as skew Gaussian, exponentially modified Gaussian, symmetric normal inverse-Gaussian, and Student's t-distribution.

J.1. Extension for Skew Gaussian

Re-expressing the update (82) in terms of μ , α , Σ (the same derivation as Khan et al. (2018)), we obtain the following update:

$$\begin{aligned}\Sigma_{t+1}^{-1} &= \Sigma_t^{-1} - 2\frac{\beta}{1-m}\nabla_{\Sigma_t}\mathcal{L} + \frac{m}{1-m}(\Sigma_t^{-1} - \Sigma_{t-1}^{-1}) \\ \mu_{t+1} &= \mu_t + \frac{\beta}{1-m}\Sigma_{t+1}\left(\frac{1}{1-c^2}\nabla_{\mu_t}\mathcal{L} - \frac{c}{1-c^2}\nabla_{\alpha_t}\mathcal{L}\right) + \frac{m}{1-m}\Sigma_{t+1}\Sigma_{t-1}^{-1}(\mu_t - \mu_{t-1}) \\ \alpha_{t+1} &= \alpha_t + \frac{\beta}{1-m}\Sigma_{t+1}\left(\frac{1}{1-c^2}\nabla_{\alpha_t}\mathcal{L} - \frac{c}{1-c^2}\nabla_{\mu_t}\mathcal{L}\right) + \frac{m}{1-m}\Sigma_{t+1}\Sigma_{t-1}^{-1}(\alpha_t - \alpha_{t-1})\end{aligned}$$

where we use a skew Gaussian distribution as the variational distribution, $\nabla_{\mu_t}\mathcal{L}$, $\nabla_{\alpha_t}\mathcal{L}$, and $\nabla_{\Sigma_t}\mathcal{L}$ are defined at (53) -(55).

We make use of the same approximations as Khan et al. (2018) such as the gradient-magnitude of the Hessian approximation, the square root approximation, $\Sigma_{t-1} \approx \Sigma_t$, and a diagonal covariance structure in Σ to obtain an extension of the variational skew-Adam update. Recall that $\mathbf{g}_\alpha^{\text{entropy}}$ and $\mathbf{g}_\Sigma^{\text{entropy}}$ are defined at (51) and (52) and $c = \sqrt{\frac{2}{\pi}}$. Using the same algebra manipulation used in Khan et al. (2018), we obtain the variational Adam update with Gaussian prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})$, where $\Sigma^{-1} = \text{Diag}(Ns + \delta)$, $v = \frac{c}{(1+\alpha^T\Sigma^{-1}\alpha)}$, and $u(\mathbf{z}) = \frac{(\mathbf{z}-\mu)^T\Sigma^{-1}\alpha}{1+\alpha^T\Sigma^{-1}\alpha}$.

Skew Gaussian extension

```

1: while not converged do
2:    $\hat{\mathbf{z}} \leftarrow \mu + \sigma \circ \epsilon$ , where  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\sigma \leftarrow 1/\sqrt{Ns + \delta}$ 
3:    $\mathbf{z} \leftarrow \hat{\mathbf{z}} + |w|\alpha$ , where  $w \sim \mathcal{N}(0, 1)$ 
4:   Randomly sample a data example  $\mathcal{D}_i$ 
5:    $\mathbf{g}_\mu \leftarrow -\nabla \log p(\mathcal{D}_i|\mathbf{z})$ 
6:   option I:  $\hat{\mathbf{g}}_\alpha \leftarrow -|w|\nabla \log p(\mathcal{D}_i|\mathbf{z})$ 
7:   option II:  $\hat{\mathbf{g}}_\alpha \leftarrow -[v\nabla \log p(\mathcal{D}_i|\hat{\mathbf{z}}) + u(\mathbf{z})\nabla \log p(\mathcal{D}_i|\mathbf{z})]$ 
8:    $\mathbf{g}_\alpha \leftarrow \hat{\mathbf{g}}_\alpha - \mathbf{g}_\alpha^{\text{entropy}}/N$ 
9:    $\mathbf{g}_s \leftarrow \mathbf{g}_\mu \circ \mathbf{g}_\mu - \text{diag}(2\mathbf{g}_\Sigma^{\text{entropy}})/N + (s + \delta/N)$ 
10:   $\mathbf{m}_\mu \leftarrow \gamma_1 \mathbf{m}_\mu + (1 - \gamma_1) \left( \frac{\mathbf{g}_\mu - c\mathbf{g}_\alpha}{1-c^2} + \delta\mu/N \right)$ 
11:   $\mathbf{m}_\alpha \leftarrow \gamma_1 \mathbf{m}_\alpha + (1 - \gamma_1) \left( \frac{\mathbf{g}_\alpha - c\mathbf{g}_\mu}{1-c^2} + \delta\alpha/N \right)$ 
12:   $\mathbf{s} \leftarrow \gamma_2 \mathbf{s} + (1 - \gamma_2) \mathbf{g}_s$ 
13:   $\hat{\mathbf{m}}_\mu \leftarrow \mathbf{m}_\mu/(1 - \gamma_1^t)$ ,  $\hat{\mathbf{m}}_\alpha \leftarrow \mathbf{m}_\alpha/(1 - \gamma_1^t)$ ,  $\hat{\mathbf{s}} \leftarrow (\mathbf{s} + \delta/N)/(1 - \gamma_2^t)$ 
14:   $\mu \leftarrow \mu - \beta \hat{\mathbf{m}}_\mu/\sqrt{\hat{\mathbf{s}}}$ 
15:   $\alpha \leftarrow \alpha - \beta \hat{\mathbf{m}}_\alpha/\sqrt{\hat{\mathbf{s}}}$ 
16:   $t \leftarrow t + 1$ 
17: end while
    
```

J.2. Extension for Exponentially Modified Gaussian

Similarly, re-expressing the update (82) in terms of μ , α , Σ , we obtain the following update:

$$\begin{aligned}\Sigma_{t+1}^{-1} &= \Sigma_t^{-1} - 2\frac{\beta}{1-m}\nabla_{\Sigma_t}\mathcal{L} + \frac{m}{1-m}(\Sigma_t^{-1} - \Sigma_{t-1}^{-1}) \\ \mu_{t+1} &= \mu_t + \frac{\beta}{1-m}\Sigma_{t+1}(2\nabla_{\mu_t}\mathcal{L} - \nabla_{\alpha_t}\mathcal{L}) + \frac{m}{1-m}\Sigma_{t+1}\Sigma_{t-1}^{-1}(\mu_t - \mu_{t-1}) \\ \alpha_{t+1} &= \alpha_t + \frac{\beta}{1-m}\Sigma_{t+1}(\nabla_{\alpha_t}\mathcal{L} - \nabla_{\mu_t}\mathcal{L}) + \frac{m}{1-m}\Sigma_{t+1}\Sigma_{t-1}^{-1}(\alpha_t - \alpha_{t-1})\end{aligned}$$

where we use an exponentially modified Gaussian distribution as the variational distribution, $\nabla_{\mu_t}\mathcal{L}$, $\nabla_{\alpha_t}\mathcal{L}$, and $\nabla_{\Sigma_t}\mathcal{L}$ are defined at (65) -(67).

Likewise, we can obtain the variational Adam update with Gaussian prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})$ as shown below, where $\mathbf{g}_\alpha^{\text{entropy}}$ and $\mathbf{g}_\Sigma^{\text{entropy}}$ are defined at (63) and (64), $\Sigma^{-1} = \text{Diag}(Ns + \delta)$, $v = \frac{1}{(\alpha^T\Sigma^{-1}\alpha)}$, and $u(\mathbf{z}) = \frac{(\mathbf{z}-\mu)^T\Sigma^{-1}\alpha-1}{\alpha^T\Sigma^{-1}\alpha}$.

Exponentially Modified Gaussian extension

```

1: while not converged do
2:    $\hat{\mathbf{z}} \leftarrow \boldsymbol{\mu} + \boldsymbol{\sigma} \circ \boldsymbol{\epsilon}$ , where  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\boldsymbol{\sigma} \leftarrow 1/\sqrt{N\mathbf{s} + \delta}$ 
3:    $\mathbf{z} \leftarrow \hat{\mathbf{z}} + w\boldsymbol{\alpha}$ , where  $w \sim \text{Exp}(1)$ 
4:   Randomly sample a data example  $\mathcal{D}_i$ 
5:    $\mathbf{g}_\mu \leftarrow -\nabla \log p(\mathcal{D}_i|\mathbf{z})$ 
6:   option I:  $\hat{\mathbf{g}}_\alpha \leftarrow -w\nabla \log p(\mathcal{D}_i|\mathbf{z})$ 
7:   option II:  $\hat{\mathbf{g}}_\alpha \leftarrow -[v\nabla \log p(\mathcal{D}_i|\hat{\mathbf{z}}) + u(\mathbf{z})\nabla \log p(\mathcal{D}_i|\mathbf{z})]$ 
8:    $\mathbf{g}_\alpha \leftarrow \hat{\mathbf{g}}_\alpha - \mathbf{g}_\alpha^{\text{entropy}}/N$ 
9:    $\mathbf{g}_s \leftarrow \mathbf{g}_\mu \circ \mathbf{g}_\mu - \text{diag}(2\mathbf{g}_\Sigma^{\text{entropy}})/N + (\mathbf{s} + \delta/N)$ 
10:   $\mathbf{m}_\mu \leftarrow \gamma_1 \mathbf{m}_\mu + (1 - \gamma_1) (2\mathbf{g}_\mu - \mathbf{g}_\alpha + \delta\boldsymbol{\mu}/N)$ 
11:   $\mathbf{m}_\alpha \leftarrow \gamma_1 \mathbf{m}_\alpha + (1 - \gamma_1) (\mathbf{g}_\alpha - \mathbf{g}_\mu + \delta\boldsymbol{\alpha}/N)$ 
12:   $\mathbf{s} \leftarrow \gamma_2 \mathbf{s} + (1 - \gamma_2) \mathbf{g}_s$ 
13:   $\hat{\mathbf{m}}_\mu \leftarrow \mathbf{m}_\mu/(1 - \gamma_1^t)$ ,  $\hat{\mathbf{m}}_\alpha \leftarrow \mathbf{m}_\alpha/(1 - \gamma_1^t)$ ,  $\hat{\mathbf{s}} \leftarrow (\mathbf{s} + \delta/N)/(1 - \gamma_2^t)$ 
14:   $\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} - \beta \hat{\mathbf{m}}_\mu/\sqrt{\hat{\mathbf{s}}}$ 
15:   $\boldsymbol{\alpha} \leftarrow \boldsymbol{\alpha} - \beta \hat{\mathbf{m}}_\alpha/\sqrt{\hat{\mathbf{s}}}$ 
16:   $t \leftarrow t + 1$ 
17: end while
    
```

J.3. Extension for Student's t-distribution

Likewise, re-expressing the update (82) in terms of $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}$, we obtain the following update:

$$\begin{aligned}\boldsymbol{\Sigma}_{t+1}^{-1} &= \boldsymbol{\Sigma}_t^{-1} - 2\frac{\beta}{1-m}\nabla_{\boldsymbol{\Sigma}_t}\mathcal{L} + \frac{m}{1-m}(\boldsymbol{\Sigma}_t^{-1} - \boldsymbol{\Sigma}_{t-1}^{-1}) \\ \boldsymbol{\mu}_{t+1} &= \boldsymbol{\mu}_t + \frac{\beta}{1-m}\boldsymbol{\Sigma}_{t+1}\nabla_{\boldsymbol{\mu}_t}\mathcal{L} + \frac{m}{1-m}\boldsymbol{\Sigma}_{t+1}\boldsymbol{\Sigma}_{t-1}^{-1}(\boldsymbol{\mu}_t - \boldsymbol{\mu}_{t-1})\end{aligned}$$

where we use a Student's t-distribution with fixed $\alpha > 1$ as the variational distribution, $\nabla_{\boldsymbol{\mu}_t}\mathcal{L}$ and $\nabla_{\boldsymbol{\Sigma}_t}\mathcal{L}$ are defined at (83) - (84).

Now, we consider the following lower bound ($\mathbf{z} \in \mathcal{R}^d$).

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})} \left[\sum_{n=1}^N \underbrace{\log p(\mathcal{D}_n|\mathbf{z})}_{-f_n(\mathbf{z})} + \log \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I}) - \log q(\mathbf{z}|\boldsymbol{\lambda}) \right].$$

where

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

We use the results from Kotz & Nadarajah (2004).

$$\begin{aligned}\mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})} [-\log q(\mathbf{z}|\boldsymbol{\lambda})] &= \frac{1}{2} \log |\boldsymbol{\Sigma}| + \log \frac{(2a\pi)^{d/2}}{\Gamma(d/2)} + \log \frac{\Gamma(d/2)\Gamma(a)}{\Gamma(d/2 + a)} + (a + d/2) (\psi(a + d/2) - \psi(a)) \\ \mathbb{E}_{q(\mathbf{z}|\boldsymbol{\lambda})} [\log \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})] &= -\frac{d}{2} \log(2\pi) + \frac{d \log(\delta)}{2} - \frac{\delta}{2} \boldsymbol{\mu}^T \boldsymbol{\mu} - \frac{\delta}{2} \frac{a}{a-1} \text{Tr}(\boldsymbol{\Sigma})\end{aligned}$$

where $\psi(\cdot)$ is the digamma function.

The remaining thing is to compute the gradients about $\mathbb{E}_{q(\mathbf{z})} [f_n(\mathbf{z})]$. To compute the gradients, the reparametrization trick can be used. However, we can do better by the extended Bonnet's and Price's theorems for Student's t-distribution (Lin et al., 2019). Assuming that $f_n(\mathbf{z})$ satisfies the assumptions needed for these two theorems, we obtain the following gradient

expression:

$$\begin{aligned}\mathbf{g}_1^n &:= \nabla_{\mu} \mathbb{E}_{q(\mathbf{z})} [f_n(\mathbf{z})] = \mathbb{E}_{q(\mathbf{z})} [\nabla_{\mathbf{z}} f_n(\mathbf{z})] \approx \nabla_{\mathbf{z}} f_n(\mathbf{z}) \\ \mathbf{g}_2^n &:= 2 \nabla_{\Sigma} \mathbb{E}_{q(\mathbf{z})} [f_n(\mathbf{z})] \\ &= \mathbb{E}_{q(\mathbf{z})} [u(\mathbf{z}) \nabla_{\mathbf{z}}^2 f_n(\mathbf{z})] \approx u(\mathbf{z}) \nabla_{\mathbf{z}}^2 f_n(\mathbf{z}) \\ &= \mathbb{E}_{q(w, \mathbf{z})} [w \nabla_{\mathbf{z}}^2 f_n(\mathbf{z})] \approx w \nabla_{\mathbf{z}}^2 f_n(\mathbf{z})\end{aligned}$$

where $\mathbf{z} \in \mathcal{R}^d$ is generated from $q(\mathbf{z})$, w is generated from $q(w)$, and

$$u(\mathbf{z}) := \frac{a + \frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{z} - \boldsymbol{\mu})}{(a + d/2 - 1)}$$

The gradients of $\mathcal{L}(\boldsymbol{\lambda})$ can be expressed as

$$\nabla_{\mu} \mathcal{L}(\boldsymbol{\lambda}) = - \sum_{n=1}^N \mathbf{g}_1^n - \delta \boldsymbol{\mu} \quad (83)$$

$$\nabla_{\Sigma} \mathcal{L}(\boldsymbol{\lambda}) = -\frac{1}{2} \sum_{n=1}^N \mathbf{g}_2^n - \frac{\delta}{2} \frac{a}{a-1} \mathbf{I} + \frac{1}{2} \boldsymbol{\Sigma}^{-1} \quad (84)$$

Likewise, we can obtain the variational Adam update with Gaussian prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})$ as shown below, where $\mathbf{z} \in \mathcal{R}^d$, $\boldsymbol{\Sigma}^{-1} = \text{Diag}(N\mathbf{s} + \frac{\alpha\delta}{\alpha-1})$, and $u(\mathbf{z}) = \frac{a + \frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{z} - \boldsymbol{\mu})}{(a + d/2 - 1)}$.

Student's t ($\alpha > 1$) extension

```

1: while not converged do
2:    $\mathbf{z} \leftarrow \boldsymbol{\mu} + \boldsymbol{\sigma} \circ \boldsymbol{\epsilon}$ , where  $w \sim \mathcal{IG}(\alpha, \alpha)$ ,  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\boldsymbol{\sigma} \leftarrow \sqrt{w/(N\mathbf{s} + \frac{\alpha\delta}{\alpha-1})}$ 
3:   Randomly sample a data example  $\mathcal{D}_i$ 
4:    $\mathbf{g}_{\mu} \leftarrow -\nabla \log p(\mathcal{D}_i|\mathbf{z})$ 
5:   option I:  $\mathbf{g}_s \leftarrow w \mathbf{g}_{\mu} \circ \mathbf{g}_{\mu}$ 
6:   option II:  $\mathbf{g}_s \leftarrow u(\mathbf{z}) \mathbf{g}_{\mu} \circ \mathbf{g}_{\mu}$ 
7:    $\mathbf{m}_{\mu} \leftarrow \gamma_1 \mathbf{m}_{\mu} + (1 - \gamma_1) (\mathbf{g}_{\mu} + \delta \boldsymbol{\mu}/N)$ 
8:    $\mathbf{s} \leftarrow \gamma_2 \mathbf{s} + (1 - \gamma_2) \mathbf{g}_s$ 
9:    $\hat{\mathbf{m}}_{\mu} \leftarrow \mathbf{m}_{\mu}/(1 - \gamma_1^t)$ ,  $\hat{\mathbf{s}} \leftarrow (\mathbf{s} + \frac{\alpha\delta}{N(\alpha-1)})/(1 - \gamma_2^t)$ 
10:   $\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} - \beta \hat{\mathbf{m}}_{\mu}/\sqrt{\hat{\mathbf{s}}}$ 
11:   $t \leftarrow t + 1$ 
12: end while
    
```

J.4. Extension for Symmetric Normal Inverse-Gaussian Distribution

Re-expressing the update (82) in terms of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ (the same derivation as Khan et al. (2018)), we obtain the following update:

$$\begin{aligned}\boldsymbol{\Sigma}_{t+1}^{-1} &= \boldsymbol{\Sigma}_t^{-1} - 2 \frac{\beta}{1-m} \nabla_{\Sigma_t} \mathcal{L} + \frac{m}{1-m} (\boldsymbol{\Sigma}_t^{-1} - \boldsymbol{\Sigma}_{t-1}^{-1}) \\ \boldsymbol{\mu}_{t+1} &= \boldsymbol{\mu}_t + \frac{\beta}{1-m} \boldsymbol{\Sigma}_{t+1} \nabla_{\mu_t} \mathcal{L} + \frac{m}{1-m} \boldsymbol{\Sigma}_{t+1} \boldsymbol{\Sigma}_{t-1}^{-1} (\boldsymbol{\mu}_t - \boldsymbol{\mu}_{t-1})\end{aligned}$$

where we use a symmetric normal inverse-Gaussian distribution with fixed $\lambda > 0$ as the variational distribution, $\nabla_{\mu_t} \mathcal{L}$ and $\nabla_{\Sigma_t} \mathcal{L}$ are defined at (78)-(79).

Likewise, we can obtain the variational Adam update with Gaussian prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \delta^{-1}\mathbf{I})$ as shown below. where

$\mathbf{z} \in \mathcal{R}^d$, $\Sigma^{-1} = \text{Diag}(N\mathbf{s} + \delta(1 + \lambda^{-1}))$, and

$$u(\mathbf{z}) = \sqrt{\frac{(\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \lambda}{\lambda}} \frac{\mathcal{K}_{\frac{d-3}{2}} \left(\sqrt{\lambda \left((\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \lambda \right)} \right)}{\mathcal{K}_{\frac{d-1}{2}} \left(\sqrt{\lambda \left((\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu}) + \lambda \right)} \right)}$$

Recall that the ratio about the Bessel functions can be approximated by Eq (71) when $d \geq 2$.

Symmetric Normal Inverse-Gaussian ($\lambda > 0$) extension

- 1: **while** not converged **do**
- 2: $\mathbf{z} \leftarrow \boldsymbol{\mu} + \boldsymbol{\sigma} \circ \boldsymbol{\epsilon}$, where $w \sim \text{InvGauss}(1, \lambda)$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $\boldsymbol{\sigma} \leftarrow \sqrt{1/[w(N\mathbf{s} + \delta(1 + \lambda^{-1}))]}$
- 3: Randomly sample a data example \mathcal{D}_i
- 4: $\mathbf{g}_\mu \leftarrow -\nabla \log p(\mathcal{D}_i | \mathbf{z})$
- 5: **option I:** $\mathbf{g}_s \leftarrow w^{-1} \mathbf{g}_\mu \circ \mathbf{g}_\mu$
- 6: **option II:** $\mathbf{g}_s \leftarrow u(\mathbf{z}) \mathbf{g}_\mu \circ \mathbf{g}_\mu$
- 7: $\mathbf{m}_\mu \leftarrow \gamma_1 \mathbf{m}_\mu + (1 - \gamma_1) (\mathbf{g}_\mu + \delta \boldsymbol{\mu}/N)$
- 8: $\mathbf{s} \leftarrow \gamma_2 \mathbf{s} + (1 - \gamma_2) \mathbf{g}_s$
- 9: $\hat{\mathbf{m}}_\mu \leftarrow \mathbf{m}_\mu / (1 - \gamma_1^t)$, $\hat{\mathbf{s}} \leftarrow (\mathbf{s} + \frac{\delta(1 + \lambda^{-1})}{N}) / (1 - \gamma_2^t)$
- 10: $\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} - \beta \hat{\mathbf{m}}_\mu / \sqrt{\hat{\mathbf{s}}}$
- 11: $t \leftarrow t + 1$
- 12: **end while**