

Scalable Inference on the Soft Affiliation Graph Model for Overlapping Community Detection

Nishma Laitonjam

NISHMA.LAITONJAM@INSIGHT-CENTRE.ORG

Weipéng Huáng

WEIPENG.HUANG@INSIGHT-CENTRE.ORG

Neil J. Hurley

NEIL.HURLEY@INSIGHT-CENTRE.ORG

Insight Centre for Data Analytics, University College Dublin, Ireland

Editors: Sinno Jialin Pan and Masashi Sugiyama

1. ELBO for S-AGM

ELBO for S-AGM is given as

$$\begin{aligned} \mathcal{L}(q) &= \mathbb{E}_q[\log p(\mathbf{A}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\pi} | \boldsymbol{\eta}, \boldsymbol{\beta})] - \mathbb{E}_q[\log q(\mathbf{W}, \boldsymbol{\pi}, \boldsymbol{\alpha})] \\ &= \sum_{ij: i < j} \mathbb{E}_q[\log p(a_{ij} | w_i, w_j, \boldsymbol{\pi})] \\ &\quad + \sum_i \sum_k \mathbb{E}_q[\log p(w_{ik} | \boldsymbol{\alpha}_k)] + \sum_k \mathbb{E}_q[\log p(\boldsymbol{\pi}_k | \eta_{k0}, \eta_{k1})] + \sum_k \mathbb{E}_q[\log p(\boldsymbol{\alpha}_k | \beta_0, \beta_1)] \\ &\quad - \sum_i \sum_k \mathbb{E}_q[\log q(w_{ik} | \phi_{ik0}, \phi_{ik1})] - \sum_k \mathbb{E}_q[\log q(\boldsymbol{\pi}_k | \lambda_{k0}, \lambda_{k1})] - \sum_k \mathbb{E}_q[\log q(\boldsymbol{\alpha}_k | \tau_{k0}, \tau_{k1})] \end{aligned}$$

Here, the term

$$\mathbb{E}_q[\log p(a_{ij} | w_i, w_j, \boldsymbol{\pi})] = a_{ij} \mathbb{E}_q[\log p_{ij}] + (1 - a_{ij}) \mathbb{E}_q[\log(1 - p_{ij})],$$

does not have an analytic form.

2. Computation of Gradients for SG-VI

In this section, the gradients i.e. $\hat{g}(\bar{\phi}_{ik0})$, $\hat{g}(\bar{\phi}_{ik1})$, $\hat{g}(\bar{\lambda}_{k0})$, $\hat{g}(\bar{\lambda}_{k1})$, $\hat{g}(\tau_{k0})$ and $\hat{g}(\tau_{k1})$ given in Algorithm 1 are computed.

2.1. Computation of $\hat{g}(\bar{\phi}_{ik0})$ and $\hat{g}(\bar{\phi}_{ik1})$

For $m \in \{0, 1\}$, we have $\hat{g}(\bar{\phi}_{ikm})$ representing the unbiased estimate of the gradient of the ELBO wrt $\bar{\phi}_{ikm}$,

$$\hat{g}(\bar{\phi}_{ikm}) = \nabla_{\bar{\phi}_{ikm}} \mathcal{L}_{\phi_{ik}}(q) = \hat{g}(\phi_{ikm}) \times \frac{1}{1 + \exp(-\bar{\phi}_{ikm})}$$

Algorithm 1 VI for the S-AGM using SGD-VI at iteration t

- 1: Sample a mini-batch \mathcal{E}^t of node pairs.
- 2: **for** Each node i in \mathcal{E}^t **do**
- 3: Sample a mini-batch of nodes \mathcal{V}_i^t .
- 4: **for** $k = 1 : K$ **do** ▷ utilizing the sampled \mathcal{V}_i^t
- 5: $\bar{\phi}_{ik0}^{(t)} = \bar{\phi}_{ik0}^{(t-1)} + \rho_\phi^{(t)} \times \hat{g}(\bar{\phi}_{ik0})$
- 6: $\bar{\phi}_{ik1}^{(t)} = \bar{\phi}_{ik1}^{(t-1)} + \rho_\phi^{(t)} \times \hat{g}(\bar{\phi}_{ik1})$
- 7: **for** $k = 1 : K$ **do** ▷ utilizing the sampled \mathcal{E}^t
- 8: $\bar{\lambda}_{k0}^{(t)} = \bar{\lambda}_{k0}^{(t-1)} + \rho_\lambda^{(t)} \times \hat{g}(\bar{\lambda}_{k0})$
- 9: $\bar{\lambda}_{k1}^{(t)} = \bar{\lambda}_{k1}^{(t-1)} + \rho_\lambda^{(t)} \times \hat{g}(\bar{\lambda}_{k1})$
- 10: **for** $k = 1 : K$ **do** ▷ utilizing the sampled \mathcal{E}^t
- 11: $\tau_{k0}^{(t)} = \tau_{k0}^{(t-1)} + \rho_\tau^{(t)} \times \hat{g}(\tau_{k0})$
- 12: $\tau_{k1}^{(t)} = \tau_{k1}^{(t-1)} + \rho_\tau^{(t)} \times \hat{g}(\tau_{k1})$

where $\hat{g}(\phi_{ikm})$ is gradient of ELBO corresponding to ϕ_{ikm} i.e. corresponding gradient wrt ϕ with preconditioning matrix $G(\phi) = \text{diag}(\phi)^{-1}$ where $\phi = \{\phi_{ik0}, \phi_{ik1}\}_{i=1,\dots,N; k=1,\dots,K}$, and $\mathcal{L}_{\phi_{ik}}(q)$ is an expression which is proportional to $\mathcal{L}(q)$ and dependent on ϕ_{ik0} and ϕ_{ik1} . We have

$$\hat{g}(\phi_{ikm}) = (G^{-1}(\phi) \nabla_\phi \mathcal{L}(q))_{ikm} = \phi_{ikm} \nabla_{\phi_{ikm}} \mathcal{L}_{\phi_{ik}}(q)$$

where,

$$\begin{aligned} \mathcal{L}_{\phi_{ik}}(q) &= \sum_{i \neq j} \mathbb{E}_q[\log p(a_{ij}|w_i, w_j, \pi)] + \mathbb{E}_q[\log p(w_{ik}|\alpha_k)] - \mathbb{E}_q[\log q(w_{ik}|\phi_{ik0}, \phi_{ik1})] \\ &= \mathbb{E}_{q(w_{ik}|\phi)} \left[\sum_{i \neq j} \mathbb{E}_{q^{-w_{ik}}}[\log p(a_{ij}|w_i, w_j, \pi)] + \mathbb{E}_{q(\alpha_k|\tau)}[\log p(w_{ik}|\alpha_k)] \right] \\ &\quad - \mathbb{E}_q[\log q(w_{ik}|\phi_{ik0}, \phi_{ik1})] \\ &= \mathbb{E}_{q(w_{ik}|\phi)} \left[\sum_{i \neq j} \mathbb{E}_{q^{-w_{ik}}}[\log p(a_{ij}|w_i, w_j, \pi)] + \Psi(\tau_{k0}) - \log(\tau_{k1}) + \left(\frac{\tau_{k0}}{\tau_{k1}} - 1 \right) \log w_{ik} \right] \\ &\quad - \mathbb{E}_q[\log q(w_{ik}|\phi_{ik0}, \phi_{ik1})] \end{aligned}$$

We can rewrite it as

$$\begin{aligned} \mathcal{L}_{\phi_{ik}}(q) &= \mathbb{E}_{q(w_{ik}|\phi)}[\hat{f}(w_{ik})] - \left(\log \Gamma(\phi_{ik0} + \phi_{ik1}) - \log \Gamma(\phi_{ik0}) - \log \Gamma(\phi_{ik1}) \right. \\ &\quad \left. + (\phi_{ik0} - 1)(\Psi(\phi_{ik0}) - \Psi(\phi_{ik0} + \phi_{ik1})) + (\phi_{ik1} - 1)(\Psi(\phi_{ik1}) - \Psi(\phi_{ik0} + \phi_{ik1})) \right) \end{aligned}$$

For a noisy gradient considering mini-batch of data i.e. considering only $j \in \mathcal{V}_i$, we have

$$\hat{f}(w_{ik}) = \frac{N}{|j \in \mathcal{V}_i^t|} \sum_{j \in \mathcal{V}_i^t} \mathbb{E}_{q^{-w_{ik}}}[\log p(a_{ij}|w_i, w_j, \pi)] + \Psi(\tau_{k0}) - \log(\tau_{k1}) + \left(\frac{\tau_{k0}}{\tau_{k1}} - 1 \right) \log w_{ik}$$

2.1.1. PARTIAL DERIVATIVE OF ELBO WRT ϕ_{ikm}

$$\begin{aligned}\nabla_{\phi_{ikm}} \mathcal{L}_{\phi_{ik}}(q) &= \nabla_{\phi_{ikm}} \mathbb{E}_{q(w_{ik}|\phi)}[\hat{f}(w_{ik})] \\ &\quad - (\phi_{ikm} - 1)\Psi'(\phi_{ikm}) + (\phi_{ik0} + \phi_{ik1} - 2)\Psi'(\phi_{ik0} + \phi_{ik1})\end{aligned}\quad (1)$$

Computation of $\nabla_{\phi_{ikm}} \mathbb{E}_{q(w_{ik}|\phi)}[\hat{f}(w_{ik})]$: For a random variable $w_{ik} \sim \text{Beta}(\phi_{ik0}, \phi_{ik1})$, we could rewrite $w_{ik} = \frac{w_{ik0}}{w_{ik0} + w_{ik1}}$ where $w_{ik0} \sim \text{Gamma}(\phi_{ik0}, 1)$ and $w_{ik1} \sim \text{Gamma}(\phi_{ik1}, 1)$. And setting $w_{ik0} = \exp(\epsilon_{ik0}\sqrt{\Psi'(\phi_{ik0})} + \Psi(\phi_{ik0}))$ and $w_{ik1} = \exp(\epsilon_{ik1}\sqrt{\Psi'(\phi_{ik1})} + \Psi(\phi_{ik1}))$, we can write

$$\nabla_{\phi_{ikm}} \mathbb{E}_{q(w_{ik}|\phi)}[\hat{f}(w_{ik})] = \nabla_{\phi_{ikm}} \mathbb{E}_{q(w_{ik0}, w_{ik1}|\phi)}[\hat{f}(w_{ik0}, w_{ik1})] = g_{\phi_{ikm}}^{rep} + g_{\phi_{ikm}}^{corr}$$

where

$$\begin{aligned}g_{\phi_{ikm}}^{rep} &= \mathbb{E}_{q(w_{ik0}, w_{ik1}; \phi)} \left[\nabla_{w_{ikm}} \hat{f}(w_{ik0}, w_{ik1}) \times \nabla_{\phi_{ikm}} \exp(\epsilon_{ikm}\sqrt{\Psi'(\phi_{ikm})} + \Psi(\phi_{ikm})) \right] \\ &= \mathbb{E}_{q(w_{ik0}, w_{ik1}; \phi)} \left[\nabla_{w_{ik}} \hat{f}(w_{ik}) \times (-1)^m w_{ik} (1 - w_{ik}) \left((\log(w_{ikm}) - \Psi(\phi_{ikm})) \times \frac{\Psi''(\phi_{ikm})}{2\Psi'(\phi_{ikm})} + \Psi'(\phi_{ikm}) \right) \right]\end{aligned}$$

and

$$\begin{aligned}g_{\phi_{ikm}}^{corr} &= \mathbb{E}_{q(w_{ik0}, w_{ik1}; \phi)} \left[\hat{f}(w_{ik0}, w_{ik1}) \left\{ \nabla_{w_{ikm}} \log q(w_{ikm}; \phi_{ikm}) \left(\nabla_{\phi_{ikm}} \exp(\epsilon_{ikm}\sqrt{\Psi'(\phi_{ikm})} + \Psi(\phi_{ikm})) \right) \right. \right. \\ &\quad \left. \left. + \nabla_{\phi_{ikm}} \log q(w_{ikm}; \phi_{ikm}) + \nabla_{\phi_{ikm}} \log J(\epsilon_{ikm}; \phi_{ikm}) \right\} \right. \\ &= \mathbb{E}_{q(w_{ik0}, w_{ik1}; \phi)} \left[\hat{f}(w_{ik0}, w_{ik1}) \left\{ (\phi_{ikm} - w_{ikm}) \left((\log(w_{ikm}) - \Psi(\phi_{ikm})) \times \frac{\Psi''(\phi_{ikm})}{2\Psi'(\phi_{ikm})} + \Psi'(\phi_{ikm}) \right) \right. \right. \\ &\quad \left. \left. + \log(w_{ikm}) - \Psi(\phi_{ikm}) + \frac{\Psi''(\phi_{ikm})}{2\Psi'(\phi_{ikm})} \right\} \right]\end{aligned}$$

where

$$\begin{aligned}J(\epsilon_{ikm}; \phi_{ikm}) &= |\det \nabla_{\epsilon_{ikm}} \exp(\epsilon_{ikm}\sqrt{\Psi'(\phi_{ikm})} + \Psi(\phi_{ikm}))| \\ &= \exp(\epsilon_{ikm}\sqrt{\Psi'(\phi_{ikm})} + \Psi(\phi_{ikm})) \sqrt{\Psi'(\phi_{ikm})}\end{aligned}$$

To compute $g_{\phi_{ikm}}^{rep}$, the term $\nabla_{w_{ik}} \hat{f}(w_{ik})$ is derived as

$$\nabla_{w_{ik}} \hat{f}(w_{ik}) = \frac{N}{|j \in \mathcal{V}_i^t|} \sum_{j \in \mathcal{V}_i^t} \mathbb{E}_{q(w_{ik})} \left[h_{ij} \frac{\pi_k w_{jk}}{1 - \pi_k w_{ik} w_{jk}} \right] + \frac{\frac{\tau_{k0}}{\tau_{k1}} - 1}{w_{ik}}$$

where

$$h_{ij} = (-1)^{1-a_{ij}} \left(p_{ij}^{-1} - 1 \right)^{a_{ij}}.$$

Putting the values of $\nabla_{\phi_{ikm}} \mathbb{E}_{q(w_{ik}|\phi)}[\hat{f}(w_{ik})]$ in Equation (1), we get $\nabla_{\phi_{ikm}} \mathcal{L}_{\phi_{ik}}(q)$ which is used to compute $\hat{g}(\bar{\phi}_{ikm})$.

Here, Γ , Ψ , ψ' and Ψ'' are gamma function, digamma function, polygamma of order 1 and polygamma of order 2 respectively.

2.2. Computation of $\hat{g}(\bar{\lambda}_{k0})$ and $\hat{g}(\bar{\lambda}_{k1})$

For $m \in \{0, 1\}$, we have $\hat{g}(\bar{\lambda}_{km})$ representing the unbiased estimate of the gradient of the ELBO wrt λ_{km} which is given by

$$\hat{g}(\bar{\lambda}_{km}) = \nabla_{\bar{\lambda}_{km}} \mathcal{L}_{\lambda_k}(q) = \hat{g}(\lambda_{km}) \times \frac{1}{1 + \exp(-\bar{\lambda}_{km})}$$

where $\hat{g}(\lambda_{km})$ is gradient of ELBO corresponding to λ_{km} i.e. corresponding gradient wrt λ with preconditioning matrix $G(\lambda) = \text{diag}(\lambda)^{-1}$ where $\lambda = \{\lambda_{k0}, \lambda_{k1}\}_{k=1, \dots, K}$, and $\mathcal{L}_{\lambda_k}(q)$ is an expression which is proportional to $\mathcal{L}(q)$ and dependent on λ_{k0} and λ_{k1} . We have

$$\hat{g}(\lambda_{km}) = (G^{-1}(\lambda) \nabla_{\lambda} \mathcal{L}(q))_{km} = \lambda_{km} \nabla_{\lambda_{km}} \mathcal{L}_{\lambda_k}(q)$$

where

$$\begin{aligned} \mathcal{L}_{\lambda_k}(q) &= \sum_{ij: i < j} \mathbb{E}_q[\log p(a_{ij}|w_i, w_j, \pi)] + \mathbb{E}_q[\log p(\pi_k|\eta_{k0}, \eta_{k1})] - \mathbb{E}_q[\log q(\pi_k|\lambda_{k0}, \lambda_{k1})] \\ &= \mathbb{E}_{q(\pi_k|\lambda_{k0}, \lambda_{k1})} \left[\sum_{ij: i < j} \mathbb{E}_{q^{-\pi_k}}[\log p(a_{ij}|w_i, w_j, \pi)] + \log p(\pi_k|\eta_{k0}, \eta_{k1}) \right] \\ &\quad - \mathbb{E}_q[\log q(\pi_k|\lambda_{k0}, \lambda_{k1})] \\ &= \mathbb{E}_{q(\pi_k|\lambda_{k0}, \lambda_{k1})} \left[\sum_{ij: i < j} \mathbb{E}_{q^{-\pi_k}}[\log p(a_{ij}|w_i, w_j, \pi)] \right. \\ &\quad \left. + \log \Gamma(\eta_{k0} + \eta_{k1}) - \log \Gamma(\eta_{k0}) - \log \Gamma(\eta_{k1}) + (\eta_{k0} - 1) \log \pi_k + (\eta_{k1} - 1) \log(1 - \pi_k) \right] \\ &\quad - \mathbb{E}_q[\log q(\pi_k|\lambda_{k0}, \lambda_{k1})] \end{aligned}$$

We can be rewrite it as

$$\begin{aligned} \mathcal{L}_{\lambda_k}(q) &= \mathbb{E}_{q(\pi_k|\lambda_{k0}, \lambda_{k1})}[\hat{f}(\pi_k)] - \left(\log \Gamma(\lambda_{k0} + \lambda_{k1}) - \log \Gamma(\lambda_{k0}) - \log \Gamma(\lambda_{k1}) \right. \\ &\quad \left. + (\lambda_{k0} - 1)(\Psi(\lambda_{k0}) - \Psi(\lambda_{k0} + \lambda_{k1})) + (\lambda_{k1} - 1)(\Psi(\lambda_{k1}) - \Psi(\lambda_{k0} + \lambda_{k1})) \right) \end{aligned}$$

For a noisy gradient considering mini-batch of data i.e. considering only $(i, j) \in \mathcal{E}^t$, we have

$$\begin{aligned} \hat{f}(\pi_k) &= s(\mathcal{E}^t) \sum_{(i,j) \in \mathcal{E}^t} \mathbb{E}_{q^{-\pi_k}}[\log p(a_{ij}|w_i, w_j, \pi)] \\ &\quad + \log \Gamma(\eta_{k0} + \eta_{k1}) - \log \Gamma(\eta_{k0}) - \log \Gamma(\eta_{k1}) + (\eta_{k0} - 1) \log \pi_k + (\eta_{k1} - 1) \log(1 - \pi_k) \end{aligned}$$

Here, $s(\mathcal{E}^t)$ is scale value for stratified sampling.

2.2.1. PARTIAL DERIVATIVE OF ELBO WRT λ_{km}

$$\begin{aligned} \nabla_{\lambda_{km}} \mathcal{L}_{\lambda_k}(q) &= \nabla_{\lambda_{km}} \mathbb{E}_{q(\pi_k|\lambda)}[\hat{f}(\pi_k)] \\ &\quad - (\lambda_{km} - 1)\Psi'(\lambda_{km}) + (\lambda_{k0} + \lambda_{k1} - 2)\Psi'(\lambda_{k0} + \lambda_{k1}) \end{aligned} \tag{2}$$

Computation of $\nabla_{\lambda_{km}} \text{E}_{q(\pi_k|\lambda_{k0},\lambda_{k1})}[\hat{f}(\pi_k)]$: For a random variable $\pi_k \sim \text{Beta}(\lambda_{k0}, \lambda_{k1})$, we could rewrite $\pi_k = \frac{\pi_{k0}}{\pi_{k0} + \pi_{k1}}$ where $\pi_{k0} \sim \text{Gamma}(\lambda_{k0}, 1)$ and $\pi_{k1} \sim \text{Gamma}(\lambda_{k1}, 1)$. And setting $\pi_{k0} = \exp(\epsilon_{k0}\sqrt{\Psi'(\lambda_{k0})} + \Psi(\lambda_{k0}))$ and $\pi_{k1} = \exp(\epsilon_{k1}\sqrt{\Psi'(\lambda_{k1})} + \Psi(\lambda_{k1}))$, we can write

$$\nabla_{\lambda_{km}} \text{E}_{q(\pi_k|\lambda)}[\hat{f}(\pi_k)] = \nabla_{\lambda_{km}} \text{E}_{q(\pi_{k0}, \pi_{k1}|\lambda)}[\hat{f}(\pi_{k0}, \pi_{k1})] = g_{\lambda_{km}}^{rep} + g_{\lambda_{km}}^{corr}$$

where

$$\begin{aligned} g_{\lambda_{km}}^{rep} &= \text{E}_{q_{\pi_{km}}(\pi_{k0}, \pi_{k1}; \lambda)} \left[\nabla_{\pi_{km}} \hat{f}(\pi_{k0}, \pi_{k1}) \times \nabla_{\lambda_{km}} \exp(\epsilon_{km}\sqrt{\Psi'(\lambda_{km})} + \Psi(\lambda_{km})) \right] \\ &= \text{E}_{q_{\pi_{km}}(\pi_{k0}, \pi_{k1}; \lambda)} \left[\nabla_{\pi_k} \hat{f}(\pi_k) \times (-1)^m \pi_k (1 - \pi_k) \left((\log(\pi_{km}) - \Psi(\lambda_{km})) \times \frac{\Psi''(\lambda_{km})}{2\Psi'(\lambda_{km})} + \Psi'(\lambda_{km}) \right) \right] \end{aligned}$$

and

$$\begin{aligned} g_{\lambda_{km}}^{corr} &= \text{E}_{q_{\pi_{km}}(\pi_{k0}, \pi_{k1}; \lambda)} \left[\hat{f}(\pi_{k0}, \pi_{k1}) \left\{ \nabla_{\pi_{km}} \log q(\pi_{km}; \lambda_{km}) \left(\nabla_{\lambda_{km}} \exp(\epsilon_{km}\sqrt{\Psi'(\lambda_{km})} + \Psi(\lambda_{km})) \right) \right. \right. \\ &\quad \left. \left. + \nabla_{\lambda_{km}} \log q(\pi_{km}; \lambda_{km}) + \nabla_{\lambda_{km}} \log J(\epsilon_{km}; \lambda_{km}) \right\} \right. \\ &= \text{E}_{q_{\pi_{km}}(\pi_{k0}, \pi_{k1}; \lambda)} \left[\hat{f}(\pi_{k0}, \pi_{k1}) \left\{ (\lambda_{km} - \pi_{km}) \left((\log(\pi_{km}) - \Psi(\lambda_{km})) \times \frac{\Psi''(\lambda_{km})}{2\Psi'(\lambda_{km})} + \Psi'(\lambda_{km}) \right) \right. \right. \\ &\quad \left. \left. + \log(\pi_{km}) - \Psi(\lambda_{km}) + \frac{\Psi''(\lambda_{km})}{2\Psi'(\lambda_{km})} \right\} \right] \end{aligned}$$

where

$$\begin{aligned} J(\epsilon_{km}; \lambda_{km}) &= |\det \nabla_{\epsilon_{km}} \exp(\epsilon_{km}\sqrt{\Psi'(\lambda_{km})} + \Psi(\lambda_{km}))| \\ &= \exp(\epsilon_{km}\sqrt{\Psi'(\lambda_{km})} + \Psi(\lambda_{km})) \sqrt{\Psi'(\lambda_{km})} \end{aligned}$$

To compute $g_{\lambda_{km}}^{rep}$, the term $\nabla_{\pi_k} \hat{f}(\pi_k)$ is derived as

$$\nabla_{\pi_k} \hat{f}(\pi_k) = s(\mathcal{E}^t) \sum_{(i,j) \in \mathcal{E}^t} \text{E}_{q^{-\pi_k}} \left[h_{ij} \frac{w_{ik} w_{jk}}{1 - \pi_k w_{ik} w_{jk}} \right] + \frac{\eta_{k0} - 1}{\pi_k} - \frac{\eta_{k1} - 1}{1 - \pi_k}$$

Putting the values of $\nabla_{\lambda_{km}} \text{E}_{q(\pi_k|\lambda_{k0},\lambda_{k1})}[\hat{f}(\pi_k)]$ in Equation (2), we get $\nabla_{\lambda_{km}} \mathcal{L}_{\lambda_k}(q)$ which is used to compute $\hat{g}(\bar{\lambda}_{km})$.

2.3. Computation of $\hat{g}(\tau_{k0})$ and $\hat{g}(\tau_{k1})$

For $m \in \{0, 1\}$, we have $\hat{g}(\tau_{km})$ representing the unbiased estimate of natural gradient of the ELBO wrt τ_{km} which is given by

$$\hat{g}(\tau_{km}) = (G^{-1}(\tau_{k0}, \tau_{k1}) \nabla_{\tau_k} \mathcal{L}_{\tau_k}(q))_m$$

which is natural gradient corresponding to τ_{km} i.e. corresponding gradient of the ELBO wrt τ_k with Fisher Information matrix as preconditioning matrix $G(\tau_{k0}, \tau_{k1})$, and $\mathcal{L}_{\tau_k}(q)$ is proportional to $\mathcal{L}(q)$ and dependent on τ_{k0} and τ_{k1} .

$$\begin{aligned}\mathcal{L}_{\tau_k}(q) &= \sum_i \mathbb{E}_q[\log p(w_{ik}|\alpha_k)] + \mathbb{E}_q[\log p(\alpha_k|\beta_0, \beta_1)] - \mathbb{E}_q[\log q(\alpha_k|\tau_{k0}, \tau_{k1})] \\ &\propto (N + \beta_0 - \tau_{k0})(\Psi(\tau_{k0}) - \log(\tau_{k1})) \\ &\quad - (\beta_1 - \sum_i (\Psi(\phi_{ik0}) - \Psi(\phi_{ik0} + \phi_{ik1})) - \tau_{k1}) \frac{\tau_{k0}}{\tau_{k1}} + \left(\log \Gamma(\tau_{k0}) - \tau_{k0} \log(\tau_{k1}) \right)\end{aligned}$$

2.3.1. PARTIAL DERIVATIVE OF ELBO WRT τ_{k0} AND τ_{km}

For a noisy gradient considering mini-batch of data i.e. considering only $i \in \mathcal{E}^t$, we have

$$\begin{aligned}\nabla_{\tau_{k0}} \mathcal{L}_{\tau_k}(q) &= \begin{bmatrix} \Psi'(\tau_{k0}) & \frac{-1}{\tau_{k1}} \end{bmatrix} \left[\beta_1 - \frac{N}{|\mathcal{E}^t|} \sum_{i \in \mathcal{E}^t} (\Psi(\phi_{ik0}) - \Psi(\phi_{ik0} + \phi_{ik1})) - \tau_{k1} \right] \\ \nabla_{\tau_{k1}} \mathcal{L}_{\tau_k}(q) &= \begin{bmatrix} \frac{-1}{\tau_{k1}} & \frac{\tau_{k0}}{\tau_{k1}^2} \end{bmatrix} \left[\beta_1 - \frac{N}{|\mathcal{E}^t|} \sum_{i \in \mathcal{E}^t} (\Psi(\phi_{ik0}) - \Psi(\phi_{ik0} + \phi_{ik1})) - \tau_{k1} \right]\end{aligned}$$

Now,

$$\begin{aligned}\nabla_{\tau_k} \mathcal{L}_{\tau_k}(q) &= \begin{bmatrix} \nabla_{\tau_{k0}} \mathcal{L}_{\tau_k}(q) \\ \nabla_{\tau_{k1}} \mathcal{L}_{\tau_k}(q) \end{bmatrix} \\ &= \begin{bmatrix} \Psi'(\tau_{k0}) & \frac{-1}{\tau_{k1}} \\ \frac{-1}{\tau_{k1}} & \frac{\tau_{k0}}{\tau_{k1}^2} \end{bmatrix} \left[\beta_1 - \frac{N}{|\mathcal{E}^t|} \sum_{i \in \mathcal{E}^t} (\Psi(\phi_{ik0}) - \Psi(\phi_{ik0} + \phi_{ik1})) - \tau_{k1} \right]\end{aligned}$$

2.3.2. NOISY NATURAL GRADIENTS, $\hat{g}(\tau_{k0})$ AND $\hat{g}(\tau_{k1})$

Using the Fisher information matrix for $\text{Gamma}(\tau_{k0}, \tau_{k1})$, i.e. $G(\tau_{k0}, \tau_{k1}) = \begin{bmatrix} \Psi'(\tau_{k0}) & \frac{-1}{\tau_{k1}} \\ \frac{-1}{\tau_{k1}} & \frac{\tau_{k0}}{\tau_{k1}^2} \end{bmatrix}$, the natural gradients are given as

$$\begin{aligned}\begin{bmatrix} \hat{g}(\tau_{k0}) \\ \hat{g}(\tau_{k1}) \end{bmatrix} &= G^{-1}(\tau_{k0}, \tau_{k1}) \nabla_{\tau_k} \mathcal{L}_{\tau_k}(q) \\ &= \left[\beta_1 - \frac{N}{|\mathcal{E}^t|} \sum_{i \in \mathcal{E}^t} (\Psi(\phi_{ik0}) - \Psi(\phi_{ik0} + \phi_{ik1})) - \tau_{k1} \right]\end{aligned}$$