

---

# Supplementary material for “Scalable Nonparametric Bayesian Multilevel Clustering”

---

**Viet Huynh** <sup>†</sup>

<sup>†</sup> Center for Pattern Recognition and Data Analytics (PRaDA)  
Deakin University, Australia

**Dinh Phung** <sup>†</sup>

**Svetha Venkatesh** <sup>†</sup>

**XuanLong Nguyen**

Department of Statistics  
University of Michigan, Ann Arbor, USA

**Matt Hoffman**

Adobe Research  
Adobe Systems, Inc.

**Hung Hai Bui**

Adobe Research  
Adobe Systems, Inc.

## Abstract

This document presents supplementary material to complement the manuscript entitled “*Scalable Nonparametric Bayesian Multilevel Clustering*”. The two first sections provide details on generalized Dirichlet distributions and some properties of exponential family. The following section presents derivations for variational and stochastic variational Bayes updates for the models.

## 1 Generalized Dirichlet distribution

Generalized Dirichlet distribution were originally introduced by (Connor & Mosiman, 1969)) and later developed with Bayesian analysis by (Wong, 1998). In this section, we introduce formal definition for this distribution and conjugate property which is useful for our derivation in the following sections.

**Definition 1.** (Original definition (Connor & Mosiman, 1969)) The generalized Dirichlet distribution (GD) is a generalization of the Dirichlet distribution with a more general covariance structure and almost twice the number of parameters. The density probability function with random vector  $\theta = (\theta_1, \dots, \theta_K)$  and  $\theta_K = 1 - \sum_{i=1}^{K-1} \theta_i$

$$p(\theta \mid a_1, \dots, a_{K-1}, b_1, \dots, b_{K-1}) = \frac{\theta_k^{b_{K-1}-1} \prod_{i=1}^{K-1} \left[ \theta_i^{a_i-1} \left( \sum_{j=i}^K \theta_j \right)^{b_{i-1}-(a_i+b_i)} \right]}{\prod_{i=1}^{K-1} B(a_i, b_i)}$$

The above distribution can be constructed as follows. If  $z_i$ 's are i.i.d. random variables from Beta distributions, i.e.  $z_i \sim \text{Beta}(a_i, b_i)$  and  $z_1 = x_1, \dots, z_i = \theta_i / \left(1 - \sum_{j=1}^{i-1} \theta_j\right)$ . We denote  $\theta \sim \text{GD}(a_1, \dots, a_{k-1}, b_1, \dots, b_{k-1})$ .

**Lemma 2.** (Conjugate prior) Suppose  $x \sim \text{GD}(a_1, \dots, a_{k-1}, b_1, \dots, b_{k-1})$  and  $y \mid x \sim \text{Mult}(x)$  then the posterior distribution is

$$x \mid y \sim \text{GD}(\tilde{a}_1, \dots, \tilde{a}_{k-1}, \tilde{b}_1, \dots, \tilde{b}_{k-1})$$

where  $\tilde{a}_i = a_i + y_i$  and  $\tilde{b}_i = b_i + \sum_{j=i+1}^k y_j$ .

*Proof.* See Wong (1998, Lemma 1). □

## 2 Some properties of Exponential family

**Proposition 3.** (Expectation of log likelihood) Let  $p(\theta \mid \eta)$  be a distribution which we call likelihood function and  $q(\theta \mid \lambda)$  be a variational distribution used to approximate  $p(\theta \mid \eta)$ . Both distributions are supposed to belong to the same exponential family form (but different (hyper)parameters), i.e.

$$p(\theta \mid \eta) \propto \exp(\langle \eta, T(\theta) \rangle - B(\eta)) \quad q(\theta \mid \lambda) \propto \exp(\langle \lambda, T(\theta) \rangle - B(\lambda))$$

Then

$$\mathbb{E}_{q(\theta|\lambda)} [\ln p(\theta \mid \eta)] = \left\langle \eta, \frac{\partial B(\lambda)}{\partial \lambda} \right\rangle - B(\eta)$$

*Proof.* We have

$$\mathbb{E}_{q(\theta|\lambda)} [\ln p(\theta \mid \eta)] = \mathbb{E}_{q(\theta|\lambda)} [\langle \eta, T(\theta) \rangle - B(\eta)]$$

However, recall the fact that  $\mathbb{E}[T(\theta)]$  is the derivative of log partition function of  $q(\theta \mid \lambda)$ , i.e.  $\mathbb{E}[T(\theta)] = \frac{\partial B(\lambda)}{\partial \lambda}$ , we have

$$\mathbb{E}_{q(\theta|\lambda)} [\ln q(\theta \mid \lambda)] = \left\langle \eta, \frac{\partial B(\lambda)}{\partial \lambda} \right\rangle - B(\eta)$$

□

**Proposition 4.** (Partial derivative of entropy) Let  $q(\theta \mid \lambda)$  be a distribution belong to exponential family, i.e.  $q(\theta \mid \lambda) \propto \exp(\langle \lambda, T(\theta) \rangle - B(\lambda))$ , then

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{q(\theta|\lambda)} [\ln q(\theta \mid \lambda)] = \frac{\partial^2 B(\lambda)}{\partial \lambda \partial \lambda^\top} \cdot \lambda$$

*Proof.* We have

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{q(\theta|\lambda)} [\ln q(\theta | \lambda)] = \frac{\partial}{\partial \lambda} \mathbb{E}_{q(\theta|\lambda)} [\langle \lambda, T(\theta) \rangle - B(\lambda)]$$

However, recall the fact that  $\mathbb{E}[T(\theta)]$  is the derivative of log partition function of  $q(\theta | \lambda)$ , i.e.  $\mathbb{E}[T(\theta)] = \frac{\partial B(\lambda)}{\partial \lambda}$ , we have

$$\begin{aligned} \frac{\partial}{\partial \lambda} \mathbb{E}_{q(\theta|\lambda)} [\ln q(\theta | \lambda)] &= \frac{\partial}{\partial \lambda} \left[ \lambda^\top \frac{\partial B(\lambda)}{\partial \lambda} - B(\lambda) \right] \\ &= \frac{\partial^2 B(\lambda)}{\partial \lambda \partial \lambda^\top} \cdot \lambda + \frac{\partial B(\lambda)}{\partial \lambda} - \frac{\partial B(\lambda)}{\partial \lambda} \\ &= \frac{\partial^2 B(\lambda)}{\partial \lambda \partial \lambda^\top} \cdot \lambda \end{aligned}$$

□

**Proposition 5.** (*Partial derivative of log likelihood*) Let  $p(\theta | \eta)$  be a distribution which we call likelihood function and  $q(\theta | \lambda)$  be a variational distribution used to approximate  $p(\theta | \eta)$ . Both distributions are supposed belong to the same exponential family form (but different (hyper)parameters), i.e.

$$p(\theta | \eta) \propto \exp(\langle \eta, T(\theta) \rangle - B(\eta)) \quad q(\theta | \lambda) \propto \exp(\langle \lambda, T(\theta) \rangle - B(\lambda))$$

Then

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{q(\theta|\lambda)} [\ln p(\theta | \eta)] = \frac{\partial^2 B(\lambda)}{\partial \lambda \partial \lambda^\top} \cdot \eta$$

*Proof.* We have

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{q(\theta|\lambda)} [\ln p(\theta | \eta)] = \frac{\partial}{\partial \lambda} \mathbb{E}_{q(\theta|\lambda)} [\langle \eta, T(\theta) \rangle - B(\eta)]$$

However, recall the fact that  $\mathbb{E}[T(\theta)]$  is the derivative of log partition function of  $q(\theta | \lambda)$ , i.e.  $\mathbb{E}[T(\theta)] = \frac{\partial B(\lambda)}{\partial \lambda}$ , we have

$$\begin{aligned} \frac{\partial}{\partial \lambda} \mathbb{E}_{q(\theta|\lambda)} [\ln p(\theta | \eta)] &= \frac{\partial}{\partial \lambda} \left[ \eta^\top \frac{\partial B(\lambda)}{\partial \lambda} - B(\eta) \right] \\ &= \frac{\partial^2 B(\lambda)}{\partial \lambda \partial \lambda^\top} \cdot \eta \end{aligned}$$

□

### 3 Variational for MC2

The objective of inference problem with the model is to estimate the posterior distribution  $p(\Theta | x, w)$  where  $\Theta$  is the collection of parameter variable of the model,  $\Theta \triangleq \{\beta, \epsilon, \tau, c, z, t, \psi, \phi\}$ . However, since this posterior is intractable, in variational Bayes inference, this will be approximated with tractable distribution called variational distribution,  $q(\Theta)$ . In order to ensure that  $q(\Theta)$ , one usually use mean-field assumption which assumes all variational variables in  $\Theta$  independent. However, because of the nature of the model, two group of variables  $z_i$  and  $t_{j1}, \dots, t_{jn_j}$  can not be totally factorized. We will maintain the joint distribution of these variables in variational inference.

The variational distribution is

$$q(\Theta) = q(\beta) q(\epsilon) q(\tau) q(c) q(z, t) q(\psi) q(\phi)$$

where (in truncation setting with  $K$  level in  $\beta$  and  $T$  level in  $\tau_k$  and  $M$  level in  $\epsilon$ ). It is noticed that  $q(z, t)$  is joint distribution in which  $t_{ji}$  conditional dependent on  $z_i$ .

$$\begin{aligned}
q(\beta) &= \text{GD}(\beta \mid \lambda^\beta) & q(\epsilon) &= \text{GD}(\epsilon \mid \lambda^\epsilon) \\
q(\tau) &= \prod_{k=1}^K \text{GD}(\tau_k \mid \lambda_k^\tau) & q(c) &= \prod_{k=1}^K \prod_{t=1}^T \text{Mult}(c_{kt} \mid \mu_{kt}^c) \\
q(\psi) &= \prod_{m=1}^M q(\psi_m \mid \lambda_m^\psi) & q(\phi) &= \prod_{k=1}^K q(\phi_k \mid \lambda_k^\phi) \\
q(z, t) &= \prod_{j=1}^J \left[ q(z_j) \prod_{i=1}^{n_j} q(t_{ji} \mid z_j) \right] = \prod_{j=1}^J \left[ \text{Mult}(z_j \mid \mu_j^z) \prod_{i=1}^{n_j} \text{Mult}(t_{ji} \mid \mu_{jiz_j}^t) \right]
\end{aligned}$$

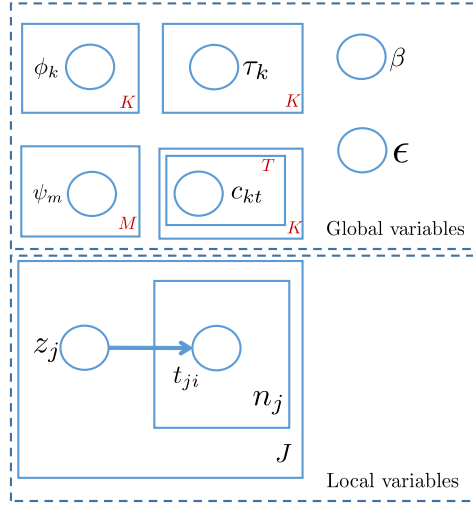


Figure 1: Variational distribution dependency.

The variational distributions is depicted in Fig. 1

- $\lambda^\beta = (\lambda_{11}^\beta, \dots, \lambda_{1(K-1)}^\beta, \lambda_{21}^\beta, \dots, \lambda_{2(K-1)}^\beta)$  is a  $2K - 2$  dimension vector;
- $\lambda^\epsilon = (\lambda_{11}^\epsilon, \dots, \lambda_{1(M-1)}^\epsilon, \lambda_{21}^\epsilon, \dots, \lambda_{2(M-1)}^\epsilon)$  is a  $2M - 2$  dimension vector;
- $\lambda^\tau = (\lambda_{11}^\tau, \dots, \lambda_{1(T-1)}^\tau, \lambda_{21}^\tau, \dots, \lambda_{2(T-1)}^\tau)$  is a  $2T - 2$  dimension vector;
- $\mu_j^z$  is a  $K$ -dimension vector;  $\mu_{kt}^c$  is a  $M$ -dimension vector;  $\mu_{jik}^t$  is a  $T$ -dimension vector.

### 3.1 Stick-breaking variable updates

Now we compute the updates for stick-breaking variables

**Update equations for  $\beta$** , we have

$$q(\beta) \propto \exp(\mathbb{E}[\ln p(x, w, \Theta)])$$

where

$$\begin{aligned}
\mathbb{E}[\ln p(x, w, \Theta)] &= \ln p(\beta) + \mathbb{E}[\ln p(z_{1:J} \mid \beta)] + \text{const} \\
&= \ln p(\beta) + \sum_{k=1}^K \sum_{j=1}^J \mu_{jk}^z \ln \beta_k + \text{const}
\end{aligned}$$

Hence,

$$q(\beta) \propto p(\beta) \exp\left(\sum_{k=1}^K \sum_{j=1}^J \mu_{jk}^z \ln \beta_k\right)$$

Since  $p(\beta | \eta)$  has the form of generalized Dirichlet distribution with parameters  $(K - 2$  dimension) $[1, \dots, 1, \gamma, \dots, \gamma]^\top$ . Therefore, using result from Lemma 2 with the prior  $p(\beta | \gamma)$  and observations  $[\sum_{j=1}^J \mu_{j1}^z, \dots, \sum_{j=1}^J \mu_{jK}^z]^\top$ , we obtain the updated distribution  $q(\beta)$  with new hyperparameter  $\lambda^\beta$  where

$$\lambda_{k1}^\beta = 1 + \sum_{j=1}^J \mu_{jk}^z \quad \lambda_{k2}^\beta = \eta + \sum_{j=1}^J \sum_{t=k+1}^K \mu_{jt}^z$$

**Update equations for  $\epsilon$ ,** The updates for  $q(\epsilon)$  can be similarly computed as follows

$$\begin{aligned} q(\epsilon) &\propto \exp(\mathbb{E}[\ln p(x, w, \Theta)]) \\ &\propto \exp(\ln p(\epsilon) + \mathbb{E}[\ln p(c | \epsilon)]) \\ &\propto \exp\left(\ln p(\epsilon) + \sum_{m=1}^M \sum_{k=1}^K \sum_{t=1}^T \epsilon_{ktm} \ln \epsilon_m\right) \\ &= p(\epsilon) \exp\left(\sum_{m=1}^M \sum_{k=1}^K \sum_{t=1}^T \epsilon_{ktm} \ln \epsilon_m\right) \end{aligned}$$

Hence, we obtain the updates

$$\lambda_{m1}^\epsilon = 1 + \sum_{k=1}^K \sum_{t=1}^T \mu_{ktm}^\epsilon \quad \lambda_{m2}^\epsilon = \gamma + \sum_{k=1}^K \sum_{t=1}^T \sum_{l=m+1}^M \mu_{ktl}^\epsilon$$

**Update equations for  $\tau_k$ ,** we have

$$\begin{aligned} \mathbb{E}[\ln p(x, w, \Theta)] &= \mathbb{E}\left[\ln p(\tau_k) + \sum_{j=1}^J (\ln p(t_j, z_j = k | \tau, \beta))\right] + \text{const} \\ &= \ln p(\tau_k) + \sum_{l=1}^l \left(\sum_{j=1}^J \mu_{jk}^z \sum_{i=1}^{n_j} \mu_{jikl}^t\right) \ln \tau_{kl} + \text{const} \end{aligned}$$

Hence,

$$q(\tau_{k,1:T}) \propto p(\tau_{k,1:T}) \exp\left(\sum_{t=1}^T \left(\sum_{j=1}^J \mu_{jk}^z \sum_{i=1}^{n_j} \mu_{jikl}^t\right) \ln \tau_{kt}\right)$$

As a consequence, we have

$$\lambda_{kl1}^\tau = 1 + \sum_{j=1}^J \mu_{jk}^z \sum_{i=1}^{n_j} \mu_{jikl}^t \quad \lambda_{kl2}^\tau = v + \sum_{j=1}^J \mu_{jk}^z \sum_{i=1}^{n_j} \sum_{r=l+1}^T \mu_{jikr}^t$$

### 3.2 Content and context atom updates

**Update equations for context atoms  $\phi_k$ ,** using the standard update for VB, we have

$$\begin{aligned}
q(\phi_k) &\propto p(\phi_k) \exp(\mathbb{E}[\ln p(x|z, \phi)]) \\
&= p(\phi_k) \exp\left(\sum_{j=1}^J \mu_{jk}^z \ln p(x_j | \phi_k)\right)
\end{aligned}$$

using conjugacy, we obtain

$$\lambda_k^\phi = \lambda_*^\phi + \sum_{j=1}^J \mu_{jk}^z [T(x_j); 1]$$

**Update equations for content atoms**  $\psi_m$ , similarly, we compute  $q(\psi_m)$  as follows

$$\begin{aligned}
q(\psi_m) &\propto p(\psi_m) \exp(\mathbb{E}[\ln p(w|c, z, t, \psi)]) \\
&= p(\psi_m) \exp\left(\sum_{j=1}^J \sum_{i=1}^{n_j} \left(\sum_{k=1}^K \vartheta_{jk} \sum_{t=1}^T \varepsilon_{ktm} \kappa_{jikt}\right) \ln p(w_{ji} | \psi_m)\right)
\end{aligned}$$

Therefore

$$\lambda_m^\psi = \lambda_*^\psi + \sum_{j=1}^J \sum_{i=1}^{n_j} \left(\sum_{k=1}^K \mu_{jk}^z \sum_{l=1}^T \mu_{ktm}^c \mu_{jikt}^t\right) [T(w_{ji}); 1]$$

### 3.3 Indicator variable updates

Now we compute variational distribution for indicators variables,  $q(z, t)$  and  $q(c)$ . The join distribution  $q(z_j, t_j)$  are factorized as follows  $q(z_j, t_j) = q(z_j) \prod_{i=1}^{n_j} q(t_{ji} | z_j)$  and compute

$$\begin{aligned}
\mu_{jikt}^t &= q(t_{ji} = l | z_j = k) \propto q(t_{ji} = l, z_j = k) \\
&\propto \exp(\mathbb{E}[\ln p(w_{ji} | z_j = k, t_{ji} = l, c, \psi)] + \mathbb{E}[\ln p(t_{ji} = l | \tau_k)]) \\
&\propto \exp\left(\sum_{m=1}^M \mu_{klm}^c \mathbb{E}[\ln p(w_{ji} | \psi_m)] + \mathbb{E}[\ln \tau_{kl}]\right)
\end{aligned}$$

This means

$$\mu_{jikt}^t \propto \exp\left(\sum_{m=1}^M \mu_{klm}^c \mathbb{E}[\ln p(w_{ji} | \psi_m)] + \mathbb{E}[\ln \tau_{kl}]\right) \triangleq \tilde{\mu}_{jikt}^t$$

In the above equation,  $\tilde{\mu}_{jikt}^t$  is the unnormalized value for  $\mu_{jikt}^t$ . This term will be used to compute  $\mu_{jk}^z$  as follows

$$\begin{aligned}
\mu_{jk}^z &= \sum_{t_{j1}, \dots, t_{jn_j}} q(z_j = k, t_j) \\
&\propto \sum_{t_{j1}, \dots, t_{jn_j}} \exp\left(\mathbb{E}[p(x_j | \phi_k)] + \mathbb{E}[\ln \beta_k] + \sum_{i=1}^{n_j} \{\mathbb{E}[\ln p(w_{ji} | z_j = k, t_{ji}, c, \psi)] + \mathbb{E}[\ln p(t_{ji} | \tau_k)]\}\right) \\
&= \exp(\mathbb{E}[p(x_j | \phi_k)] + \mathbb{E}[\ln \beta_k]) \left(\sum_{t_{j1}, \dots, t_{jn_j}} \prod_i \exp(\{\mathbb{E}[\ln p(w_{ji} | z_j = k, t_{ji}, c, \psi)] + \mathbb{E}[\ln p(t_{ji} | \tau_k)]\})\right) \\
&= \exp(\mathbb{E}[p(x_j | \phi_k)] + \mathbb{E}[\ln \beta_k]) \prod_i \left(\sum_{l=1}^T \tilde{\mu}_{jikt}^t\right) \\
&= \exp\left(\mathbb{E}[p(x_j | \phi_k)] + \mathbb{E}[\ln \beta_k] + \sum_i \ln \left(\sum_{l=1}^T \tilde{\mu}_{jikt}^t\right)\right)
\end{aligned}$$

Therefore, we obtain the update equation for  $\mu_j^z$  as follows

$$\mu_{jk}^z \propto \exp \left( \mathbb{E} [\ln p(x_j | \phi_k)] + \mathbb{E} [\ln \beta_k] + \sum_i \ln \left( \sum_{l=1}^T \tilde{\mu}_{j^t ikl}^t \right) \right)$$

Finally, we need to compute update for  $\mu_{klm}^c$ , we have

$$\begin{aligned} \mu_{klm}^c &\propto \exp (\mathbb{E} [\ln p(x, w, \Theta)]) \\ &\propto \exp (\mathbb{E} [\ln p(w | c_{kl} = m, z, t, \psi)] + \mathbb{E} [\ln \epsilon_m]) \\ &= \exp \left( \sum_{j=1}^J \mu_{jk}^z \sum_{i=1}^{n_j} \mu_{j^t ikl}^t \mathbb{E} [\ln p(x_{ji} | \psi_m)] + \mathbb{E} [\ln \epsilon_m] \right) \end{aligned}$$

which is

$$\mu_{klm}^c \propto \exp \left( \sum_{j=1}^J \mu_{jk}^z \sum_{i=1}^{n_j} \mu_{j^t ikl}^t \mathbb{E} [\ln p(x_{ji} | \psi_m)] + \mathbb{E} [\ln \epsilon_m] \right)$$

## 4 Stochastic Variational for MC2

The Evidence Lower Bound (ELBO) function for the model

$$\mathcal{F}(x, w, \Theta) \triangleq \mathbb{E} \{ \ln p(x, w, \Theta) \} - \mathbb{E} \{ \ln q(\Theta) \}$$

can be represented as  $\mathcal{F} = \sum_{j=1}^J \mathcal{F}_j = \mathbb{E} [J\mathcal{F}_j]$  where  $J$  is the number of observation groups;  $\mathcal{F}_j$  is lower bound function only related to document  $j$ -th and defined as

$$\begin{aligned} \mathcal{L}_j &= \mathbb{E} [\ln p(x_j | z_j, \phi)] + \mathbb{E} [\ln p(w_j | z_j, t_j, c, \psi)] + \mathbb{E} [\ln p(z_j | \beta)] + \mathbb{E} [\ln p(t_j | z_j, \tau)] \\ &\quad - \mathbb{E} [\ln q(z_j)] - \mathbb{E} [\ln q(t_j | z_j)] \\ &\quad + \frac{1}{J} (\mathbb{E} [\ln p(c | \epsilon)] + \mathbb{E} [\ln p(\beta)] + \mathbb{E} [\ln p(\tau)] + \mathbb{E} [\ln p(\epsilon)]) \\ &\quad + \frac{1}{J} (\mathbb{E} [\ln p(\psi | \lambda_*^\psi)] + \mathbb{E} [\ln p(\phi | \lambda_*^\phi)]) \\ &\quad - \frac{1}{J} (\mathbb{E} [\ln q(c)] + \mathbb{E} [\ln q(\beta)] + \mathbb{E} [\ln q(\tau)] + \mathbb{E} [\ln q(\epsilon)]) \\ &\quad - \frac{1}{J} (\mathbb{E} [\ln q(\psi)] + \mathbb{E} [\ln q(\phi)]) \end{aligned}$$

### 4.1 Stochastic updates for stick-breaking variables

**Update equations for  $\beta$** , using results from Proposition 4 and 5 with generalized Dirichlet distributions, we get

$$\frac{\partial}{\partial \lambda_k^\beta} \mathbb{E} [\ln q(\beta)] = \frac{\partial^2 B(\lambda_k^\beta)}{\partial \lambda_k^\beta \partial (\lambda_k^\beta)^\top} [\lambda_{k1}^\beta - 1; \lambda_{k2}^\beta - 1]$$

and

$$\frac{\partial}{\partial \lambda_k^\beta} \mathbb{E} [\ln p(\beta)] = \frac{\partial^2 B(\lambda_k^\beta)}{\partial \lambda_k^\beta \partial (\lambda_k^\beta)^\top} [0; \gamma - 1]$$

Therefore, using natural gradient, we have

$$\frac{\partial^{(\text{ng})} \mathcal{L}_j}{\partial \lambda_{k1}^\beta} = \frac{-\lambda_{k1}^\beta + 1 + J \mu_{jk}^z}{J} \quad \frac{\partial^{(\text{ng})} \mathcal{L}_j}{\partial \lambda_{k2}^\beta} = \frac{-\lambda_{k2}^\beta + \eta + J \sum_{l=k+1}^K \mu_{jl}^z}{J}$$

**Update equations for  $\epsilon$ ,** we have

$$\frac{\partial \mathcal{L}_j}{\partial \lambda_m^\epsilon} = \frac{\partial}{\partial \lambda_m^\epsilon} \left( \frac{1}{J} \mathbb{E} [\ln p(c | \epsilon)] + \frac{1}{J} \mathbb{E} [\ln p(\epsilon)] - \frac{1}{J} \mathbb{E} [\ln q(\epsilon)] \right)$$

where

$$\begin{aligned} \frac{\partial}{\partial \lambda_m^\epsilon} \mathbb{E} [\ln p(c | \epsilon)] &= \frac{\partial}{\partial \lambda_m^\epsilon} \mathbb{E} \left[ \sum_{k=1}^K \sum_{t=1}^T \sum_{r=1}^M \mu_{ktr}^c \ln p(c_{lt} = r | \epsilon) \right] \\ &= \frac{\partial}{\partial \lambda_m^\epsilon} \mathbb{E} \left[ \left\langle [\ln \epsilon_m, \ln(1 - \epsilon_m)], \left[ \sum_{k=1}^K \sum_{t=1}^T \mu_{ktm}^c, \sum_{k=1}^K \sum_{t=1}^T \sum_{r=m+1}^M \mu_{ktr}^c \right] \right\rangle \right] \\ &= \frac{\partial^2 B(\lambda_m^\epsilon)}{\partial \lambda_m^\epsilon \partial (\lambda_m^\epsilon)^\top} \left[ \sum_{k=1}^K \sum_{t=1}^T \mu_{ktm}^c, \sum_{k=1}^K \sum_{t=1}^T \sum_{r=m+1}^M \mu_{ktr}^c \right] \end{aligned}$$

Similarly

$$\frac{\partial}{\partial \lambda_m^\epsilon} \mathbb{E} [\ln q(\epsilon)] = \frac{\partial^2 B(\lambda_m^\epsilon)}{\partial \lambda_m^\epsilon \partial (\lambda_m^\epsilon)^\top} [\lambda_{m1}^\epsilon - 1, \lambda_{m2}^\epsilon - 1]$$

and

$$\frac{\partial}{\partial \lambda_m^\epsilon} \mathbb{E} [\ln p(u)] = \left\langle \frac{\partial^2 B(\varphi_m)}{\partial \varphi_m \partial \varphi_m^\top}, [0; \gamma - 1] \right\rangle$$

Therefore, using natural gradient, we have

$$\frac{\partial^{(\text{ng})} \mathcal{L}_j}{\partial \lambda_{m1}^\epsilon} = \frac{-\lambda_{m1}^\epsilon + 1 + \sum_{k=1}^K \sum_{t=1}^T \mu_{ktm}^c}{J} \quad \frac{\partial^{(\text{ng})} \mathcal{L}_j}{\partial \lambda_{m2}^\epsilon} = \frac{-\lambda_{m2}^\epsilon + \gamma + \sum_{k=1}^K \sum_{t=1}^T \sum_{r=m+1}^M \mu_{ktr}^c}{J}$$

**Update equations for  $\tau_k$ ,** we have

$$\frac{\partial \mathcal{L}_j}{\partial \lambda_{kt}^\tau} = \frac{\partial}{\partial \lambda_{kt}^\tau} \left( \mathbb{E} [\ln p(t_j | \tau, z_j)] + \frac{1}{J} \mathbb{E} [\ln p(\tau)] - \frac{1}{J} \mathbb{E} [\ln q(\tau)] \right)$$

where

$$\begin{aligned} \frac{\partial}{\partial \lambda_{kl}^\tau} \mathbb{E} [\ln p(t_j | \tau, z_j)] &= \mu_{jk}^z \sum_{i=1}^{n_j} \frac{\partial}{\partial \lambda_{kl}^\tau} \mathbb{E} [\mu_{jikl}^t \ln p(t_{ji} = l | \tau_k, z_j = k)] \\ &= \mu_{jk}^z \frac{\partial}{\partial \lambda_{kl}^\tau} \mathbb{E} \left[ \sum_{i=1}^{n_j} \mu_{jikl}^t \ln(\tau_k) \right] \\ &= \mu_{jk}^z \frac{\partial}{\partial \lambda_{kl}^\tau} \mathbb{E} \left[ \left\langle [\ln \tau_{kl}, \ln(1 - \tau_{kl})], \left[ \sum_{i=1}^{n_j} \mu_{jikl}^t, \sum_{i=1}^{n_j} \sum_{r=l+1}^T \mu_{jikr}^t \right] \right\rangle \right] \\ &= \mu_{jk}^z \frac{\partial^2 B(\lambda_{kl}^\tau)}{\partial \lambda_{kl}^\tau \partial (\lambda_{kl}^\tau)^\top} \left[ \sum_{i=1}^{n_j} \mu_{jikl}^t, \sum_{i=1}^{n_j} \sum_{r=l+1}^T \mu_{jikr}^t \right] \end{aligned}$$



Similarly

$$\frac{\partial}{\partial \lambda_{kl}^\tau} \mathbb{E} [\ln q(\tau_{kl})] = \frac{\partial^2 B(\lambda_{kl}^\tau)}{\partial \lambda_{kl}^\tau \partial (\lambda_{kl}^\tau)^\top} [\lambda_{kl1}^\tau - 1; \lambda_{kl2}^\tau - 1]$$

and

$$\frac{\partial}{\partial \lambda_{kl}^\tau} \mathbb{E} [\ln p(\tau_{kl})] = \frac{\partial^2 B(\lambda_{kl}^\tau)}{\partial \lambda_{kl}^\tau \partial (\lambda_{kl}^\tau)^\top} [0; \nu - 1]$$

Finally, using natural gradient, we have

$$\frac{\partial^{(\text{ng})} \mathcal{L}_j}{\partial \lambda_{kl1}^\tau} = \frac{-\lambda_{kl1}^\tau + 1 + J \mu_{jk}^z \sum_{i=1}^{n_j} \mu_{jikl}^t}{J} \quad \frac{\partial^{(\text{ng})} \mathcal{L}_j}{\partial \lambda_{kl2}^\tau} = \frac{-\lambda_{kl2}^\tau + \nu + J \mu_{jk}^z \sum_{i=1}^{n_j} \sum_{r=t+1}^T \mu_{jikr}^t}{J}$$

## 4.2 Stochastic updates for content and context atoms

**Update equations for content atoms**  $\phi_k$ , compute gradient of  $\mathcal{L}_j$  with respect to  $\lambda_k^\phi$  (in  $q(\phi_k | \lambda_k^\phi)$ ), we have

$$\begin{aligned} \frac{\partial \mathcal{L}_j}{\partial \lambda_k^\phi} &= \frac{\partial}{\partial \lambda_k^\phi} \left( \mathbb{E} [\ln p(x_j | z, \phi)] + \frac{1}{J} \left( \mathbb{E} [\ln p(\phi | \lambda_*^\phi)] - \mathbb{E} [\ln q(\phi)] \right) \right) \\ &= \frac{\partial}{\partial \lambda_k^\phi} \left( \mathbb{E} [\ln p(x_j | z, \phi)] + \frac{1}{J} \left( \mathbb{E} [\ln p(\phi_k | \lambda_*^\phi)] - \mathbb{E} [\ln q(\phi_k | \lambda_k^\phi)] \right) \right) \end{aligned}$$

where

$$\begin{aligned} \frac{\partial}{\partial \lambda_k^\phi} \mathbb{E} [\ln p(x_j | z, \phi)] &= \frac{\partial}{\partial \lambda_k^\phi} \sum_{k=1}^K \mu_{jk}^z \mathbb{E} [\langle T(x_j), \phi_k \rangle - A(\phi_k)] \\ &= \frac{\partial}{\partial \lambda_k^\phi} \mu_{jk}^z \left\langle [T(x_j), 1], \frac{\partial B(\lambda_k^\phi)}{\partial \lambda_k^\phi} \right\rangle \\ &= \mu_{jk}^z \left\langle [T(x_j), 1], \frac{\partial^2 \lambda_k^\phi}{\partial \lambda_k^\phi \partial (\lambda_k^\phi)^\top} \right\rangle \end{aligned}$$

here we use the fact that  $\mathbb{E}[\phi_k, -A(\phi_k)]$  is the derivative of log partition function of  $q(\phi_k | \lambda_k^\phi)$ , i.e.  $\frac{\partial B(\lambda_k^\phi)}{\partial \lambda_k^\phi}$ .

Moreover, using results from Proposition 4 and 5, we get

$$\frac{\partial}{\partial \lambda_k^\phi} \mathbb{E} [\ln p(\phi | \lambda_*^\phi)] = \left\langle \lambda_*^\phi, \frac{\partial^2 \lambda_k^\phi}{\partial \lambda_k^\phi \partial (\lambda_k^\phi)^\top} \right\rangle$$

and

$$\frac{\partial}{\partial \lambda_k^\phi} \mathbb{E} [\ln q(\phi_k | \lambda_k^\phi)] = \left\langle \lambda_k^\phi, \frac{\partial^2 \lambda_k^\phi}{\partial \lambda_k^\phi \partial (\lambda_k^\phi)^\top} \right\rangle$$

Hence, using natural gradient, we have

$$\frac{\partial^{(\text{ng})} \mathcal{L}_j}{\partial \lambda_k^\phi} = \frac{-\lambda_k^\phi + \lambda_*^\phi + J \mu_{jk}^z [T(x_j), 1]}{J}$$

**Update equations for content atoms**  $\psi_m$ , computing gradient of  $\mathcal{L}_j$  with respect to  $\lambda_m^\psi$  (in  $q(\psi_m | \lambda_m^\psi)$ ), we have

$$\begin{aligned} \frac{\partial \mathcal{L}_j}{\partial \lambda_m^\psi} &= \frac{\partial}{\partial \lambda_m^\psi} \left( \mathbb{E} [\ln p(w_j | z, t, c, \psi)] + \frac{1}{J} \left( \mathbb{E} [\ln p(\psi | \lambda_*^\psi)] - \mathbb{E} [\ln q(\psi)] \right) \right) \\ &= \frac{\partial}{\partial \lambda_m^\psi} \left( \mathbb{E} [\ln p(w_j | z, t, c, \psi)] + \frac{1}{J} \left( \mathbb{E} [\ln p(\psi_m | \lambda_*^\psi)] - \mathbb{E} [\ln q(\psi_m | \lambda_m^\psi)] \right) \right) \end{aligned}$$

where

$$\begin{aligned} \frac{\partial}{\partial \lambda_m^\psi} \mathbb{E} [\ln p(w_j | z, t, c, \psi)] &= \frac{\partial}{\partial \lambda_m^\psi} \sum_{i=1}^{n_j} \sum_{k=1}^K \mu_{jk}^z \sum_{l=1}^T \mu_{jikl}^t \sum_{m=1}^M \mu_{klm}^c \mathbb{E} [\langle T \langle w_{ji} \rangle, \psi_m \rangle - A(\psi_m)] \\ &= \frac{\partial}{\partial \lambda_m^\psi} \sum_{k=1}^K \mu_{jk}^z \sum_{l=1}^T \mu_{klm}^c \sum_{i=1}^{n_j} \mu_{jikl}^t \left\langle [T \langle w_{ji} \rangle; 1], \frac{\partial B(\lambda_m^\psi)}{\partial \lambda_m^\psi} \right\rangle \\ &= \sum_{k=1}^K \mu_{jk}^z \sum_{l=1}^T \mu_{klm}^c \sum_{i=1}^{n_j} \mu_{jikl}^t \frac{\partial^2 B(\lambda_m^\psi)}{\partial \lambda_m^\psi \partial (\lambda_m^\psi)^\top} [T \langle w_{ji} \rangle; 1] \end{aligned}$$

Here we use the fact that  $\mathbb{E}[\psi_m, -A(\psi_m)]$  is the derivative of log partition function of  $q(\psi_m | \lambda_m^\psi)$ , i.e.  $\frac{\partial B(\lambda_m^\psi)}{\partial \lambda_m^\psi}$ . Moreover, using results from Proposition 4 and 5, we get

$$\frac{\partial}{\partial \lambda_m^\psi} \mathbb{E} [\ln p(\psi_m | \lambda_*^\psi)] = \frac{\partial^2 B(\lambda_m^\psi)}{\partial \lambda_m^\psi \partial (\lambda_m^\psi)^\top} \lambda_*^\psi$$

and

$$\frac{\partial}{\partial \lambda_m^\psi} \mathbb{E} [\ln q(\psi_m | \lambda_m)] = \frac{\partial^2 B(\lambda_m^\psi)}{\partial \lambda_m^\psi \partial (\lambda_m^\psi)^\top} \lambda_m^\psi$$

Hence, using natural gradient, we have

$$\frac{\partial^{(\text{ng})} \mathcal{L}_j}{\partial \lambda_m^\psi} = \frac{-\lambda_m^\psi + \lambda_*^\psi + J \sum_{k=1}^K \mu_{jk}^z \sum_{l=1}^T \mu_{klm}^c \sum_{i=1}^{n_j} \mu_{jikl}^t [T \langle w_{ji} \rangle; 1]}{J}$$

### 4.3 Stochastic updates for global indicator variables

Since updating  $\mu_{kt}^c$  involving computation from all data groups  $j$ 's, we will use “lazy update” with stochastic gradient. The gradient of ELBO function over  $\mu_{kt}^c$  is needed to compute. However, the mean-parameterization of  $q(c_{kt} | \mu_{kt}^c)$  with constraints  $\sum_m \mu_{ktm}^c = 1$  and  $0 \leq \mu_{ktm}^c \leq 1$  is difficult to compute gradient. We therefore prefer to work with a minimal natural parameterization in exponential family form as follows

$$q(c_{kt} | \lambda_{kt}^c) = \exp(\langle \lambda_{kt}^c, T(c_{kt}) \rangle - B(\lambda_{kt}^c))$$

where  $\lambda_{kt}^c = [\lambda_{kt1}^c, \dots, \lambda_{kt(M-1)}^c]^\top$ ;  $T(c) = [\delta(c-1), \dots, \delta(c-M+1)]^\top$  and  $B(\lambda_{kt}^c) = 1 + \sum_{m=1}^{M-1} \exp(\lambda_{ktm}^c)$ . The relationship between these groups of parameters is

$$\begin{aligned} \mu_{ktm}^c &= \frac{\exp(\lambda_{ktm}^c)}{1 + \sum_{m=1}^{M-1} \exp(\lambda_{ktm}^c)}, m = 1, \dots, M-1 \\ \text{and } \mu_{kM}^c &= \frac{1}{1 + \sum_{m=1}^{M-1} \exp(\lambda_{ktm}^c)} \end{aligned}$$

Now we can compute the gradient over new parameters  $\lambda_{kt}^c$ :

$$\frac{\partial^{(\text{ng})} \mathcal{L}_j}{\partial \lambda_{klm}^c} = \frac{-\lambda_{klm}^c + \mathbb{E} \left[ \ln \frac{\epsilon_m}{\epsilon_M} \right]}{J} + (a_{klm} - a_{klM}) \quad (1)$$

where  $a_{klm} = \mu_{jk}^z \sum_{i=1}^{n_j} \mu_{jikt}^t \mathbb{E} [\ln p(w_{ji} | \psi_m)]$ , for all  $m$ .

## References

- Connor, R. J., & Mosiman, J. E. 1969. Concepts of independence for proportions with a generalization of the Dirichlet distribution. *Journal of the American Statistical Association*, **64**, 194–206.
- Wong, T.-T. 1998. Generalized Dirichlet distribution in Bayesian analysis. *Applied Mathematics and Computation*, **97**, 165–181.