

VIBP derivations

Joey

2021

1 Introduction

Existing Variational Indian Buffet Process methods do not update the hyper-parameters, such as α , σ_A , and σ_n (Doshi et al., 2009). Also, it's using the mean field approximation rather than the stochastic variant. This note is for deriving the equations myself.

2 Derivations

Directly start with the infinite linear gaussian model.

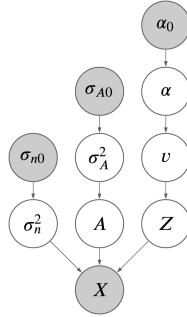


Figure 1: (a) True model $p(\mathbf{W})$

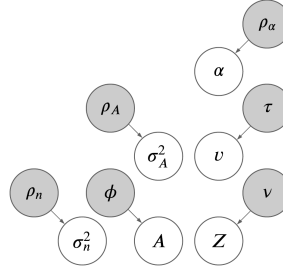


Figure 2: (b) Variational approximation $q(\mathbf{W})$

Denote the set of hidden variables in IBP by $\mathbf{W} = \{\mathbf{v}, \mathbf{Z}, \mathbf{A}, \alpha, \sigma_A^2, \sigma_n^2\}$ (σ_A^2 and σ_n^2 are precision, for simpler update formula later. Also, consider them as a single term own, not a squared term) and the set of hyper-parameters by $\boldsymbol{\theta} = \{\alpha_0, \sigma_{A0}, \sigma_{n0}\}$. The true log posterior

$$\log p(\mathbf{W}|\mathbf{X}, \boldsymbol{\theta}) = \log p(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta}) - \log p(\mathbf{X}|\boldsymbol{\theta}) \quad (1)$$

is difficult due to the intractability of computing the log marginal probability $\log p(\mathbf{X}|\boldsymbol{\theta}) = \log \int \log p(\mathbf{X}, \mathbf{W}|\boldsymbol{\theta}) \mathbf{W}$

Mean field variational methods approximate the true posterior with a variational distribution $q_{\Phi}(\mathbf{W})$, where Φ denotes the set of parameters $\{\phi, \tau, \nu, \rho_n, \rho_A, \rho_{\alpha}\}$.

Since σ_A^2 , σ_n^2 , and α are all positive, I set them to follow a Gamma distribution. This will prevent the troublesome term if using a normal distribution. To be clear, the hyper-parameter of Gamma distribution, for example, α_0 , is a vector includes $\{\alpha_s, \alpha_r\}$, where s stands for shape and r stands for rate. The model for p is the full stick-breaking construction for IBP:

$$\begin{aligned}
\alpha &\sim \text{Gamma}(\alpha_s, \alpha_r) \\
\sigma_A^2 &\sim \text{Gamma}(\sigma_{As}, \sigma_{Ar}) \\
\sigma_n^2 &\sim \text{Gamma}(\sigma_{ns}, \sigma_{nr}) \\
v_k &\sim \text{Beta}(\alpha, 1) \\
\pi_k &= \prod_{i=1}^k v_i \\
z_{nk} &\sim \text{Bernoulli}(\pi_k) \\
\mathbf{A}_{k\cdot} &\sim \text{Normal}(0, \sigma_A^2 \mathbf{I}) \quad \sigma_A^2 \text{ is precision} \\
\mathbf{X}_n &\sim \text{Normal}(\mathbf{Z}_n \mathbf{A}, \sigma_n^2 \mathbf{I}) \quad \sigma_n^2 \text{ is precision}
\end{aligned} \tag{2}$$

The joint is

$$\begin{aligned}
p(\mathbf{W}, \mathbf{X} | \theta) &= \prod_{k=1}^{\infty} \left(p(v_k | \alpha) p(\mathbf{A}_{k\cdot} | \sigma_A^2 \mathbf{I}) \prod_{n=1}^N p(z_{nk} | v_k) \right) \\
&\quad \prod_{n=1}^N p(\mathbf{X}_n | \mathbf{Z}_n, \mathbf{A}, \sigma_n^2 \mathbf{I}) p(\alpha | \alpha_0) p(\sigma_A^2 | \sigma_{A0}) p(\sigma_n^2 | \sigma_{n0}) \tag{3}
\end{aligned}$$

We use variational approximation. A truncated stick-breaking process is used here (Blei and Jordan, 2004).

Our mean field variational distribution is:

$$q(\mathbf{W}) = q_{\rho_{\alpha}}(\alpha) q_{\rho_A}(\sigma_A^2) q_{\rho_n}(\sigma_n^2) q_{\tau}(\mathbf{v}) q_{\theta}(\mathbf{A}) q_{\nu}(\mathbf{Z}) \tag{4}$$

where

$$\begin{aligned}
q_{\rho_{\alpha}}(\alpha) &= \text{Gamma}(\alpha; \rho_{\alpha s}, \rho_{\alpha r}) \\
q_{\rho_A}(\sigma_A^2) &= \text{Gamma}(\sigma_A^2; \rho_{As}, \rho_{Ar}) \\
q_{\rho_n}(\sigma_n^2) &= \text{Gamma}(\sigma_n^2; \rho_{ns}, \rho_{nr}) \\
q_{\tau_k}(v_k) &= \text{Beta}(v_k; \tau_{k1}, \tau_{k2}) \\
q_{\phi_k}(\mathbf{A}_{k\cdot}) &= \text{Normal}(\mathbf{A}_{k\cdot}; \bar{\phi}_k, \Phi_k) \quad \Phi_k \text{ is variance} \\
q_{\nu_{nk}}(z_{nk}) &= \text{Bernoulli}(z_{nk}; \nu_{nk})
\end{aligned} \tag{5}$$

Inference involves optimising parameters to minimise the KL divergence $D(q||p)$, or equivalently to maximise the ELBO.

2.1 Lower Bound on the Marginal Likelihood

Start from the KL divergence (Kingma and Welling, 2013; Odaibo, 2019),

$$\begin{aligned}
D_{KL}(q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) || p(\mathbf{W}|\mathbf{X}, \theta)) &\geq 0 \\
&= - \int q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \log \left(\frac{p(\mathbf{W}|\mathbf{X}, \theta)}{q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} \right) d\mathbf{W} \\
&= - \int q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \log \left(\frac{p(\mathbf{X}|\mathbf{W}, \theta)p(\mathbf{W}|\theta)}{q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)p(\mathbf{X}|\theta)} \right) d\mathbf{W} \\
&= - \int q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \left[\log \left(\frac{p(\mathbf{X}|\mathbf{W}, \theta)p(\mathbf{W}|\theta)}{q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} \right) - \log p(\mathbf{X}|\theta) \right] d\mathbf{W} \\
&= - \int q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \log \left(\frac{p(\mathbf{X}|\mathbf{W}, \theta)p(\mathbf{W}|\theta)}{q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} \right) d\mathbf{W} + \int q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \log p(\mathbf{X}|\theta) d\mathbf{W} \\
&= - \int q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \log \left(\frac{p(\mathbf{X}|\mathbf{W}, \theta)p(\mathbf{W}|\theta)}{q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} \right) d\mathbf{W} + \log p(\mathbf{X}|\theta) \int q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) d\mathbf{W} \\
&= - \int q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \log \left(\frac{p(\mathbf{X}|\mathbf{W}, \theta)p(\mathbf{W}|\theta)}{q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} \right) d\mathbf{W} + \log p(\mathbf{X}|\theta) \geq 0
\end{aligned} \tag{6}$$

reordering and applying rules of logarithms,

$$\begin{aligned}
\log p(\mathbf{X}|\theta) &\geq \int q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \left[\log p(\mathbf{X}|\mathbf{W}, \theta) + \log p(\mathbf{W}|\theta) - \log q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \right] d\mathbf{W} \\
&= \mathbb{E}_{\mathbf{W} \sim q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} \left[\log p(\mathbf{X}|\mathbf{W}, \theta) + \log p(\mathbf{W}|\theta) - \log q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \right] \\
&= \mathbb{E}_{\mathbf{W} \sim q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} \left[\log p(\mathbf{X}, \mathbf{W}|\theta) - \log q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \right] \\
&= \mathbb{E}_{\mathbf{W} \sim q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} [\log p(\mathbf{X}, \mathbf{W}|\theta)] + H[q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)]
\end{aligned} \tag{7}$$

where H is the entropy of distribution q . The term on the RHS is usually called variational lower bound on the marginal likelihood or the Evidence Lower Bound (ELBO). From Equation 6, we also have

$$\begin{aligned}
\log p(\mathbf{X}|\theta) &\geq \int q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \log \left(\frac{p(\mathbf{X}|\mathbf{W}, \theta)p(\mathbf{W}|\theta)}{q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} \right) d\mathbf{W} \\
&= \int q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \log \left(\frac{p(\mathbf{W}|\theta)}{q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} \right) d\mathbf{W} + \int q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) \log p(\mathbf{X}|\mathbf{W}, \theta) d\mathbf{W} \\
&= -D_{KL}(q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) || \log p(\mathbf{W}|\theta)) + \mathbb{E}_{\mathbf{W} \sim q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} [\log p(\mathbf{X}|\mathbf{W}, \theta)]
\end{aligned} \tag{8}$$

Therefore,

$$\begin{aligned}
&\arg \min_{\tau, \theta, \nu, \rho_{\alpha}, \rho_A, \rho_n} D_{KL}(q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) || p(\mathbf{W}|\mathbf{X}, \theta)) \\
&= \arg \max_{\tau, \theta, \nu, \rho_{\alpha}, \rho_A, \rho_n} -D_{KL}(q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta) || \log p(\mathbf{W}|\theta)) + \mathbb{E}_{\mathbf{W} \sim q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} [\log p(\mathbf{X}|\mathbf{W}, \theta)]
\end{aligned} \tag{9}$$

2.2 Decomposing

Continue from Equation 7, the RHS can be decomposed as

$$\begin{aligned}
&\mathbb{E}_{\mathbf{W} \sim q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)} [\log p(\mathbf{X}, \mathbf{W}|\theta)] + H[q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)] \\
&= \mathbb{E}_{\alpha \sim q} [\log p(\alpha|\alpha_0)] + \mathbb{E}_{\sigma_A^2 \sim q} [\log p(\sigma_A^2|\sigma_{A0})] + \mathbb{E}_{\sigma_n^2 \sim q} [\log p(\sigma_n^2|\sigma_{n0})] \\
&\quad + \sum_{k=1}^K \mathbb{E}_{v, \alpha \sim q} [\log p(v_k|\alpha)] + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}_{v, \mathbf{Z} \sim q} [\log p(Z_{nk}|\mathbf{v})] \\
&\quad + \sum_{k=1}^K \mathbb{E}_{\mathbf{A}, \sigma_A^2 \sim q} [\log p(\mathbf{A}_k|\sigma_A^2 I)] + \sum_{n=1}^N \mathbb{E}_{\mathbf{Z}, \mathbf{A}, \sigma_n^2 \sim q} [\log p(\mathbf{X}_n|\mathbf{Z}_n, \mathbf{A}, \sigma_n^2 I)] \\
&\quad + H[q_{\Phi}(\mathbf{W}|\mathbf{X}, \theta)]
\end{aligned} \tag{10}$$

To make it shorter, we will omit $\sim q$ in expectation term if possible from now on.

1-3. The first three terms are similar, since they all follow Gamma distribution. So here I derive use a random variable x as an example for α , σ_A^2 , and σ_n^2 .

Suppose X has $\Gamma(\alpha, \beta)$ distribution and we'd like to get the expectation of $Y = \log(X)$. Because β is a rate parameter, it will shift by $-\log(\beta)$. Therefore we start with the case $\beta = 1$. The probability element of X is

$$f_X(x) = \frac{1}{\Gamma(\alpha)} x^{\alpha-1} e^{-x} \tag{11}$$

where $\Gamma(\alpha) = \int_0^{\infty} x^\alpha e^{-x} \frac{dx}{x}$ is a constant. Substitute $x = e^y$, and $\frac{dx}{x} = dy$

$$f_Y(y) = \frac{1}{\Gamma(\alpha)e^{\alpha y - e^x}} dy \quad (12)$$

which is a differentiable function of α . Gives

$$\frac{d}{d\alpha} e^{\alpha y - e^x} dy = y \Gamma(\alpha) f_Y(y) \quad (13)$$

So,

$$\begin{aligned} \mathbb{E}(Y) &= \int y f_Y(y) = \frac{1}{\Gamma(\alpha)} \int \frac{d}{d\alpha} e^{\alpha y - e^x} dy \\ &= \frac{1}{\Gamma(\alpha)} \frac{d}{d\alpha} \int e^{\alpha y - e^x} dy \\ &= \frac{1}{\Gamma(\alpha)} \frac{d}{d\alpha} \Gamma(\alpha) \\ &= \frac{d}{d\alpha} \log \Gamma(\alpha) \\ &= \psi(\alpha) \end{aligned} \quad (14)$$

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

Reintroducing β

$$\mathbb{E}(\log(X)) = -\log(\beta) + \psi(\alpha) \quad (15)$$

Therefore, the expectation we want is

$$\begin{aligned} \mathbb{E}_{x \sim q(\alpha_q, \beta_q)} \left[\log p(x | \alpha_p, \beta_p) \right] &= \mathbb{E}_{\sim q} \left[\log \left(\frac{\beta_p^{\alpha_p}}{\Gamma(\alpha_p)} x^{\alpha_p - 1} e^{-\beta_p x} \right) \right] \\ &= \mathbb{E} \left[\frac{\beta_p^{\alpha_p}}{\Gamma(\alpha_p)} \right] + (\alpha_p - 1) \mathbb{E}[\log x] - \beta_p \mathbb{E}[x] \\ &= \frac{\beta_p^{\alpha_p}}{\Gamma(\alpha_p)} + (\alpha_p - 1)(-\log(\beta_q) + \psi(\alpha_q)) - \beta_p \frac{\alpha_q}{\beta_q} \end{aligned} \quad (16)$$

As to α , σ_A^2 , and σ_n^2 ,

$$\begin{aligned} \mathbb{E}_\alpha [\log p(\alpha | \alpha_0)] &= \frac{\alpha_r^{\alpha_s}}{\Gamma(\alpha_s)} + (\alpha_s - 1)(-\log(\rho_{\alpha r}) + \psi(\rho_{\alpha s})) - \alpha_s \frac{\rho_{\alpha s}}{\rho_{\alpha r}} \\ \mathbb{E}_{\sigma_A^2} [\log p(\sigma_A^2 | \sigma_{A0})] &= \frac{\sigma_{Ar}^{\sigma_{As}}}{\Gamma(\sigma_{As})} + (\sigma_{As} - 1)(-\log(\rho_{Ar}) + \psi(\rho_{As})) - \sigma_{As} \frac{\rho_{As}}{\rho_{Ar}} \\ \mathbb{E}_{\sigma_n^2} [\log p(\sigma_n^2 | \sigma_{n0})] &= \frac{\sigma_{nr}^{\sigma_{ns}}}{\Gamma(\sigma_{ns})} + (\sigma_{ns} - 1)(-\log(\rho_{nr}) + \psi(\rho_{ns})) - \sigma_{ns} \frac{\rho_{ns}}{\rho_{nr}} \end{aligned} \quad (17)$$

4. According to the derivation in the previous part, the forth term becomes

$$\begin{aligned}
\mathbb{E}_{\mathbf{v}, \alpha}[\log p(v_k | \alpha)] &= \mathbb{E}_{\mathbf{v}, \alpha}[\log(\alpha v_k^{\alpha-1})] \\
&= \mathbb{E}_{\alpha}[\log \alpha] + \mathbb{E}_{\alpha}[(\alpha - 1)] \mathbb{E}_{\mathbf{v}}[\log(v_k)] \\
&= (-\log(\rho_{\alpha r}) + \psi(\rho_{\alpha s})) + \left(\frac{\rho_{\alpha s}}{\rho_{\alpha r}} - 1\right) (\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2}))
\end{aligned} \tag{18}$$

5. The fifth term becomes

$$\begin{aligned}
\mathbb{E}_{\mathbf{v}, \mathbf{Z}}[\log p(Z_{nk} | \mathbf{v})] &= \mathbb{E}_{\mathbf{v}, \mathbf{Z}}[\log p(z_{nk} = 1 | \mathbf{v})^{z_{nk}} p(z_{nk} = 0 | \mathbf{v})^{1-z_{nk}}] \\
&= \mathbb{E}_{\mathbf{v}, \mathbf{Z}}[z_{nk} \log p(z_{nk} = 1 | \mathbf{v}) + (1 - z_{nk}) \log p(z_{nk} = 0 | \mathbf{v})] \\
&= \mathbb{E}_{\mathbf{Z}} z_{nk} \mathbb{E}_{\mathbf{v}} \log p(z_{nk} = 1 | \mathbf{v}) + \mathbb{E}_{\mathbf{Z}} (1 - z_{nk}) \mathbb{E}_{\mathbf{v}} \log p(z_{nk} = 0 | \mathbf{v}) \\
&= \nu_{nk} \mathbb{E}_{\mathbf{v}} \left[\log \prod_{m=1}^k v_m \right] + (1 - \nu_{nk}) \mathbb{E}_{\mathbf{v}} \left[\log \left(1 - \prod_{m=1}^k v_m \right) \right] \\
&= \nu_{nk} \left(\sum_{m=1}^k (\psi(\tau_{m1}) - \psi(\tau_{m1} + \tau_{m2})) \right) + (1 - \nu_{nk}) \mathbb{E}_{\mathbf{v}} \left[\log \left(1 - \prod_{m=1}^k v_m \right) \right]
\end{aligned} \tag{19}$$

The last expectation can be bounded using an auxiliary multinomial approximation $q_k = (q_{k1}, q_{k2}, \dots, q_{kk})$. The detailed derivation is in Doshi et al., 2009 Section 5.1. Only list the final expression here

$$\begin{aligned}
\mathbb{E}_{\mathbf{v}} \left[\log \left(1 - \prod_{k=1}^k v_m \right) \right] &\geq \left(\sum_{m=1}^k q_{km} \psi(\tau_{m2}) \right) + \left(\sum_{m=1}^{k-1} \left(\sum_{n=m+1}^k q_{kn} \right) \psi(\tau_{m1}) \right) \\
&\quad - \left(\sum_{m=1}^k \left(\sum_{n=m}^k q_{kn} \right) \psi(\tau_{m1} + \tau_{m2}) \right) - \sum_{m=1}^k q_{km} \log q_{km}
\end{aligned} \tag{20}$$

where

$$q_{ki} \propto \exp \left(\psi(\tau_{i2}) + \sum_{m=1}^{i-1} \psi(\tau_{m1}) - \sum_{m=1}^i \psi(\tau_{m1} + \tau_{m2}) \right) \tag{21}$$

6. First, the expectation of the inverse of a $Gamma(\alpha, \beta)$ distributed random variable is required.

$$\begin{aligned}
\mathbb{E}_{x \sim p(\alpha, \beta)} \left[\frac{1}{X} \right] &= \int \frac{1}{x} \beta \frac{(\beta x)^{\alpha-1}}{\Gamma(\alpha)} e^{-\beta x} dx \\
&= \beta \frac{\Gamma(\alpha-1)}{\Gamma(\alpha)} \int \frac{(\beta x)^{\alpha-2}}{\Gamma(\alpha-1)} e^{-\beta x} dx \\
&= \beta \frac{(\alpha-1)!}{(\alpha-2)!} \cdot 1 \\
&= \frac{\beta}{\alpha-1}
\end{aligned} \tag{22}$$

The integral is 1 because it's the density of $\text{Gamma}(\alpha-1, \beta)$.

Also, according to the multivariate quadratic form

$$\mathbb{E}[\mathbf{A}^T \mathbf{\Lambda} \mathbf{A}] = \text{tr}[\mathbf{\Lambda} \mathbf{\Sigma}] + \mu^T \mathbf{\Lambda} \mu \tag{23}$$

where μ and $\mathbf{\Sigma}$ are the expected value and covariance matrix of vector \mathbf{A} . Therefore, the sixth term becomes

$$\begin{aligned}
\mathbb{E}_{\mathbf{A}, \sigma_A^2} [\log p(\mathbf{A}_k | 0, \sigma_A^2 I)] &= \mathbb{E}_{\mathbf{A}, \sigma_A^2} \left[\log \left(\frac{\sigma_A^2 D/2}{(2\pi)^{D/2}} e^{-\frac{1}{2} \mathbf{A}_k^T \sigma_A^2 \mathbf{A}_k} \right) \right] \\
&= \frac{D}{2} \mathbb{E}_{\sigma_A^2} [\log(\frac{\sigma_A^2}{2\pi})] - \frac{1}{2} \mathbb{E}_{\sigma_A^2} [\sigma_A^2] \mathbb{E}_{\mathbf{A}} [\mathbf{A}_k^T \mathbf{A}_k] \\
&= \frac{D}{2} \left(\mathbb{E}_{\sigma_A^2} \log(\sigma_A^2) - \log(2\pi) \right) - \frac{1}{2} \mathbb{E}_{\sigma_A^2} [\sigma_A^2] \mathbb{E}_{\mathbf{A}} [\mathbf{A}_k^T \mathbf{A}_k] \\
&= \frac{D}{2} \left((-\log(\rho_{Ar}) + \psi(\rho_{As}) - \log(2\pi)) - \frac{\rho_{As}}{2\rho_{Ar}} \left(\text{tr}(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T \right) \right)
\end{aligned} \tag{24}$$

7. The seventh term, which is the likelihood, becomes

$$\begin{aligned}
& \mathbb{E}_{\mathbf{Z}, \mathbf{A}, \sigma_n^2} [\log(\mathbf{X}_{n\cdot} | \mathbf{Z}_{n\cdot}, \mathbf{A}, \sigma_n^2)] \\
&= \mathbb{E}_{\mathbf{Z}, \mathbf{A}, \sigma_n^2} \left[\log \left(\left(\frac{\sigma_n^2}{2\pi} \right)^{D/2} e^{-\frac{1}{2}(\mathbf{X}_{n\cdot} - \mathbf{Z}_{n\cdot} \mathbf{A}) \sigma_n^2 (\mathbf{X}_{n\cdot} - \mathbf{Z}_{n\cdot} \mathbf{A})^T} \right) \right] \\
&= \mathbb{E}_{\sigma_n^2} [\log(\frac{\sigma_n^2}{2\pi})^{D/2}] - \frac{1}{2} \mathbb{E}_{\sigma_n^2} [\sigma_n^2] \mathbb{E}_{\mathbf{Z}, \mathbf{A}} [(\mathbf{X}_{n\cdot} - \mathbf{Z}_{n\cdot} \mathbf{A})(\mathbf{X}_{n\cdot} - \mathbf{Z}_{n\cdot} \mathbf{A})^T] \\
&= \frac{D}{2} \left((-\log(\rho_{nr}) + \psi(\rho_{ns})) - \log(2\pi) \right) \\
&\quad - \frac{1}{2} \mathbb{E}_{\sigma_n^2} [\sigma_n^2] \left(\mathbf{X}_{n\cdot} \mathbf{X}_{n\cdot}^T - 2 \mathbb{E}_{\mathbf{Z}} [\mathbf{Z}_{n\cdot}] \mathbb{E}_{\mathbf{A}} [\mathbf{A}] \mathbf{X}_{n\cdot}^T + \mathbb{E}_{\mathbf{Z}, \mathbf{A}} [\mathbf{Z}_{n\cdot} \mathbf{A} \mathbf{A}^T \mathbf{Z}_{n\cdot}^T] \right) \\
&= \frac{D}{2} \left((-\log(\rho_{nr}) + \psi(\rho_{ns})) - \log(2\pi) \right) \\
&\quad - \left(\frac{\rho_{ns}}{2\rho_{nr}} \right) \left(\mathbf{X}_{n\cdot} \mathbf{X}_{n\cdot}^T - 2 \sum_{k=1}^K \nu_{nk} \bar{\phi}_k \mathbf{X}_{n\cdot}^T + 2 \sum_{k < k'} \nu_{nk} \nu_{nk'} \bar{\phi}_k \bar{\phi}_{k'}^T + \sum_{k=1}^K \nu_{nk} (tr(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T) \right)
\end{aligned} \tag{25}$$

where the last expectation is derived by

$$\begin{aligned}
\mathbb{E}_{\mathbf{Z}, \mathbf{A}} [\mathbf{Z}_{n\cdot} \mathbf{A} \mathbf{A}^T \mathbf{Z}_{n\cdot}^T] &= \mathbb{E}_{\mathbf{Z}, \mathbf{A}} \left[\left(\sum_{k=1}^K z_{nk} \mathbf{A}_{k\cdot} \right) \left(\sum_{k=1}^K z_{nk} \mathbf{A}_{k\cdot} \right)^T \right] \\
&= \mathbb{E}_{\mathbf{Z}, \mathbf{A}} \left[\sum_{d=1}^D \left(\sum_{k=1}^K z_{nk}^2 A_{kd}^2 + 2 \sum_{k, k': k' \neq k} z_{nk} z_{nk'} \mathbf{A}_{kd} \mathbf{A}_{k'd} \right)^T \right] \\
&= \sum_{k=1}^K \nu_{nk} (tr(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T) + 2 \sum_{k < k'} \nu_{nk} \nu_{nk'} \bar{\phi}_k \bar{\phi}_{k'}^T
\end{aligned} \tag{26}$$

8. The last term entropy

$$\begin{aligned}
H[q] &= -\mathbb{E}_q \log \left[q_{\rho_\alpha}(\alpha) q_{\rho_A}(\sigma_A^2) q_{\rho_n}(\sigma_n^2) \prod_{k=1}^K q_{\tau_k}(\mathbf{v}_k) \prod_{k=1}^K q_{\theta_k}(\mathbf{A}_{k\cdot}) \prod_{k=1}^K \prod_{n=1}^N q_{\nu_{nk}}(z_{nk}) \right] \\
&= \mathbb{E}_\alpha (-\log q_{\rho_\alpha}(\alpha)) + \mathbb{E}_{\sigma_A^2} (-\log q_{\rho_A}(\sigma_A^2)) + \mathbb{E}_{\sigma_n^2} (-\log q_{\rho_n}(\sigma_n^2)) \\
&\quad + \sum_{k=1}^K \mathbb{E}_{\mathbf{v}} (-\log q_{\tau_k}(\mathbf{v}_k)) + \sum_{k=1}^K \mathbb{E}_{\mathbf{A}} (-\log q_{\theta_k}(\mathbf{A}_{k\cdot})) + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}_{\mathbf{Z}} (-\log q_{\nu_{nk}}(z_{nk}))
\end{aligned} \tag{27}$$

The first three term are still in the same form. So again, as an example,

$$\begin{aligned}
\mathbb{E}_{x \sim \Gamma(\alpha, \beta)}(-\log(q(x))) &= \mathbb{E}_x \left[-\log \left(\frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} \right) \right] \\
&= \log \left(\frac{\Gamma(\alpha)}{\beta^\alpha} \right) - (\alpha - 1) \mathbb{E}[\log x] + \beta \mathbb{E}[x] \\
&= \log \left(\frac{\Gamma(\alpha)}{\beta^\alpha} \right) - (\alpha - 1)(-\log(\beta) + \psi(\alpha)) + \alpha
\end{aligned} \tag{28}$$

Therefore, we can write down each term explicitly as below

$$\begin{aligned}
\mathbb{E}_\alpha(-\log q_{\rho_\alpha}(\alpha)) &= \log \left(\frac{\Gamma(\rho_{\alpha s})}{\rho_{\alpha r}^{\rho_{\alpha s}}} \right) - (\rho_{\alpha s} - 1)(-\log(\rho_{\alpha r}) + \psi(\rho_{\alpha s})) + \rho_{\alpha s} \\
\mathbb{E}_{\sigma_A^2}(-\log q_{\rho_A}(\sigma_A^2)) &= \log \left(\frac{\Gamma(\rho_{As})}{\rho_{Ar}^{\rho_{As}}} \right) - (\rho_{As} - 1)(-\log(\rho_{Ar}) + \psi(\rho_{As})) + \rho_{As} \\
\mathbb{E}_{\sigma_n^2}(-\log q_{\rho_n}(\sigma_n^2)) &= \log \left(\frac{\Gamma(\rho_{ns})}{\rho_{nr}^{\rho_{ns}}} \right) - (\rho_{ns} - 1)(-\log(\rho_{nr}) + \psi(\rho_{ns})) + \rho_{ns} \\
\mathbb{E}_{\mathbf{v}}(-\log q_{\tau_k}(v_k)) &= \log \left(\frac{\Gamma(\tau_{k1})\Gamma(\tau_{k2})}{\Gamma(\tau_{k1} + \tau_{k2})} \right) \\
&\quad - (\tau_{k1} - 1)\psi(\tau_{k1}) - (\tau_{k2} - 1)\psi(\tau_{k2}) + (\tau_{k1} + \tau_{k2} - 2)\psi(\tau_{k1} + \tau_{k2}) \\
\mathbb{E}_{\mathbf{A}}(-\log q_{\phi_k}(\mathbf{A}_{k\cdot})) &= \frac{1}{2} \log((2\pi e)^D |\Phi_k|) \\
\mathbb{E}_{\mathbf{Z}}(-\log q_{\nu_{nk}}(z_{nk})) &= -\nu_{nk} \log \nu_{nk} - (1 - \nu_{nk}) \log(1 - \nu_{nk})
\end{aligned} \tag{29}$$

Put them together gives the variational lower bound.

2.3 Parameter Updates

Since they are in the exponential family, we use an alternate formula to update the parameters (Beal, 2003)

For any variational parameters ξ_i that correspond to W_i , the optimal ξ_i are the solution to

$$\log q_{\xi_i}(W_i) = \mathbb{E}_{\mathbf{W}_{-i}}[\log p(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta})] + c \tag{30}$$

1. Updates the feature weight matrix \mathbf{A}

$$\begin{aligned}
\log q_{\phi_k}(\mathbf{A}_k) &= \mathbb{E}_{\mathbf{A}_{-k}, \mathbf{Z}, \sigma_A^2, \sigma_n^2} [\log p(\mathbf{W}, \mathbf{X} | \boldsymbol{\theta})] + c \\
&= \mathbb{E}_{\mathbf{A}_{-k}, \mathbf{Z}, \sigma_A^2, \sigma_n^2} \left[\log p(\mathbf{A}_{k\cdot} | \sigma_A^2) + \sum_{n=1}^N \log p(\mathbf{X}_{n\cdot} | \mathbf{Z}_{n\cdot}, \mathbf{A}, \sigma_n^2) \right] + c \\
&= \mathbb{E}_{\mathbf{A}_{-k}, \mathbf{Z}, \sigma_A^2, \sigma_n^2} \left[\log \left(\left(\frac{\sigma_A^2}{2\pi} \right)^{D/2} e^{-\frac{\mathbf{A}_{k\cdot} \sigma_A^2 \mathbf{A}_{k\cdot}^T}{2}} \right) \right. \\
&\quad \left. + \sum_{n=1}^N \log \left(\left(\frac{\sigma_n^2}{2\pi} \right)^{D/2} e^{-\frac{1}{2} (\mathbf{X}_n - \mathbf{Z}_{n\cdot} \mathbf{A}) \sigma_n^2 (\mathbf{X}_n - \mathbf{Z}_{n\cdot} \mathbf{A})^T} \right) \right] + c \\
&= \frac{D}{2} \mathbb{E}_{\sigma_A^2} [\sigma_A^2] \mathbb{E}_{\mathbf{A}_{-k}} (\mathbf{A}_{k\cdot} \mathbf{A}_{k\cdot}^T) \\
&\quad - \frac{1}{2} \mathbb{E}_{\sigma_n^2} [\sigma_n^2] \sum_{n=1}^N \mathbb{E}_{\mathbf{A}_{-k}, \mathbf{Z}} [(\mathbf{X}_n - \mathbf{Z}_{n\cdot} \mathbf{A}) (\mathbf{X}_n - \mathbf{Z}_{n\cdot} \mathbf{A})^T] + c \\
&= \frac{D}{2} \mathbb{E}_{\sigma_A^2} [\sigma_A^2] (\mathbf{A}_{k\cdot} \mathbf{A}_{k\cdot}^T) \\
&\quad - \frac{1}{2} \mathbb{E}_{\sigma_n^2} [\sigma_n^2] \sum_{n=1}^N \mathbb{E}_{\mathbf{A}_{-k}, \mathbf{Z}} [(\mathbf{X}_n - \mathbf{Z}_{n(-k)} \mathbf{A}_{-k} - \mathbf{Z}_{nk} \mathbf{A}_k) (\mathbf{X}_n - \mathbf{Z}_{n(-k)} \mathbf{A}_{-k} - \mathbf{Z}_{nk} \mathbf{A}_k)^T] + c \\
&= \frac{D\rho_{As}}{2\rho_{Ar}} (\mathbf{A}_{k\cdot} \mathbf{A}_{k\cdot}^T) \\
&\quad - \frac{\rho_{ns}}{2\rho_{nr}} \sum_{n=1}^N \mathbb{E}_{\mathbf{A}_{-k}, \mathbf{Z}} \left[-2\mathbf{Z}_{nk} \mathbf{A}_{k\cdot} (\mathbf{X}_{n\cdot} - \mathbf{Z}_{n(-k)} \mathbf{A}_{-k})^T + \mathbf{Z}_{nk} \mathbf{A}_{k\cdot} \mathbf{A}_{k\cdot}^T \mathbf{Z}_{nk}^T \right] + c \\
&= \frac{D\rho_{As}}{2\rho_{Ar}} (\mathbf{A}_{k\cdot} \mathbf{A}_{k\cdot}^T) \\
&\quad - \frac{\rho_{ns}}{2\rho_{nr}} \sum_{n=1}^N \left(-2\nu_{nk} \mathbf{A}_{k\cdot} \left(\mathbf{X}_{n\cdot} - \sum_{l:l \neq k} \nu_{nl} \cdot \bar{\boldsymbol{\phi}}_l \right) + \mathbb{E}[\mathbf{Z}_{nk}^2] (\mathbf{A}_{k\cdot} \mathbf{A}_{k\cdot}^T) \right) + c \\
&= \frac{D\rho_{As}}{2\rho_{Ar}} (\mathbf{A}_{k\cdot} \mathbf{A}_{k\cdot}^T) \\
&\quad - \frac{\rho_{ns}}{2\rho_{nr}} \sum_{n=1}^N \left(-2\nu_{nk} \mathbf{A}_{k\cdot} \left(\mathbf{X}_{n\cdot} - \sum_{l:l \neq k} \nu_{nl} \cdot \bar{\boldsymbol{\phi}}_l \right)^T + \nu_{nk} (\mathbf{A}_{k\cdot} \mathbf{A}_{k\cdot}^T) \right) + c \\
&= \frac{1}{2} \left(\mathbf{A}_{k\cdot} \left(\frac{D\rho_{As}}{\rho_{Ar}} - \frac{\rho_{ns} \sum_{n=1}^N \nu_{nk}}{\rho_{nr}} \right) \mathbf{A}_{k\cdot}^T - 2\mathbf{A}_{k\cdot} \frac{\rho_{ns} \sum_{n=1}^N \nu_{nk}}{\rho_{nr}} \left(\mathbf{X}_{n\cdot} - \sum_{l:l \neq k} \nu_{nl} \cdot \bar{\boldsymbol{\phi}}_l \right)^T \right) + c
\end{aligned} \tag{31}$$

q_ϕ is supposed to be a matrix normal. Therefore, it gives the updates

$$\begin{aligned}
\bar{\phi}_k &= \left[\frac{\rho_{ns} \sum_{n=1}^N \nu_{nk}}{\rho_{nr}} \left(\mathbf{X}_{n\cdot} - \sum_{l:l \neq k} \nu_{nl} \cdot \bar{\phi}_l \right) \right] \left(\frac{D\rho_{As}}{\rho_{Ar}} - \frac{\rho_{ns} \sum_{n=1}^N \nu_{nk}}{\rho_{nr}} \right)^{-1} \\
\Phi_k &= \left(\frac{D\rho_{As}}{\rho_{Ar}} - \frac{\rho_{ns} \sum_{n=1}^N \nu_{nk}}{\rho_{nr}} \right)^{-1} I
\end{aligned} \tag{32}$$

2. Updates the feature matrix \mathbf{Z} at position nk

$$\begin{aligned}
\log q_{\nu_{nk}}(z_{nk}) &= \mathbb{E}_{\mathbf{v}, \mathbf{A}, \mathbf{Z}_{-nk}, \sigma_n^2} [\log p(\mathbf{W}, \mathbf{X} | \boldsymbol{\theta})] + c \\
&= \mathbb{E}_{\mathbf{v}, \mathbf{A}, \mathbf{Z}_{-nk}, \sigma_n^2} [\log p(z_{nk} | v_k) + \log p(\mathbf{X}_n | \mathbf{Z}_{n\cdot}, \mathbf{A}, \sigma_n^2 I)] + c
\end{aligned} \tag{33}$$

where the first term

$$\begin{aligned}
\mathbb{E}_{\mathbf{v}, \mathbf{Z}_{-nk}} [\log p(z_{nk} | v_k)] &= \mathbb{E}_{\mathbf{Z}_{-nk}} z_{nk} \mathbb{E}_{\mathbf{v}} \log p(z_{nk} = 1 | \mathbf{v}) + \mathbb{E}_{\mathbf{Z}_{-nk}} (1 - z_{nk}) \mathbb{E}_{\mathbf{v}} \log p(z_{nk} = 0 | \mathbf{v}) \\
&= z_{nk} \left(\sum_{m=1}^k (\psi(\tau_{m1}) - \psi(\tau_{m1} + \tau_{m2})) \right) + (1 - z_{nk}) \mathbb{E}_{\mathbf{v}} \left[\log \left(1 - \prod_{m=1}^k v_m \right) \right]
\end{aligned} \tag{34}$$

The remaining expectation is approximated following equation 19. The second term

$$\begin{aligned}
\mathbb{E}_{\mathbf{A}, \mathbf{Z}_{-nk}, \sigma_n^2} [\log p(\mathbf{X}_n | \mathbf{Z}_{n\cdot}, \mathbf{A}, \sigma_n^2 I)] &= -\frac{1}{2} \mathbb{E}_{\sigma_n^2} [\sigma_n^2] \left(-2 \mathbb{E}_{\mathbf{Z}_{-nk}} [\mathbf{Z}_{n\cdot}] \mathbb{E}_{\mathbf{A}} [\mathbf{A}] \mathbf{X}_{n\cdot}^T + \mathbb{E}_{\mathbf{Z}_{-nk}, \mathbf{A}} [\mathbf{Z}_{n\cdot} \mathbf{A} \mathbf{A}^T \mathbf{Z}_{n\cdot}^T] \right) + c \\
&= -\left(\frac{\rho_{ns}}{2\rho_{nr}} \right) \left(-2 z_{nk} \bar{\phi}_k \mathbf{X}_{n\cdot}^T + \mathbb{E}_{\mathbf{Z}_{-nk}, \mathbf{A}} [\mathbf{Z}_{n\cdot} \mathbf{A} \mathbf{A}^T \mathbf{Z}_{n\cdot}^T] \right) + c
\end{aligned} \tag{35}$$

The last term is decomposed as follow ($z_{nk}^2 = z_{nk}$, since z is binary)

$$\begin{aligned}
\mathbb{E}_{\mathbf{Z}_{-nk}, \mathbf{A}} [\mathbf{Z}_{n\cdot} \mathbf{A} \mathbf{A}^T \mathbf{Z}_{n\cdot}^T] &= \mathbb{E}_{\mathbf{Z}_{-nk}, \mathbf{A}} \left[\left(\sum_{k=1}^K z_{nk} \mathbf{A}_{k\cdot} \right) \left(\sum_{k=1}^K z_{nk} \mathbf{A}_{k\cdot} \right)^T \right] \\
&= \mathbb{E}_{\mathbf{Z}_{-nk}, \mathbf{A}} \left[\sum_{d=1}^D \left(\sum_{k=1}^K z_{nk}^2 A_{kd}^2 + 2 \sum_{k, k': k' \neq k} z_{nk} z_{nk'} A_{kd} A_{k'd} \right) \right]^T \\
&= z_{nk}^2 (tr(\Phi_k) + \bar{\phi}_k \bar{\phi}_k^T) + 2 z_{nk} \bar{\phi}_k \sum_{k' \neq k} \nu_{nk'} \bar{\phi}_{k'}^T \\
&= z_{nk} (tr(\Phi_k) + \bar{\phi}_k \bar{\phi}_k^T) + 2 z_{nk} \bar{\phi}_k \sum_{k' \neq k} \nu_{nk'} \bar{\phi}_{k'}^T
\end{aligned} \tag{36}$$

Therefore,

$$\begin{aligned}\mathbb{E}_{\mathbf{v}, \mathbf{Z}_{-nk}}[\log p(z_{nk}|v_k)] &= z_{nk} \left(\sum_{m=1}^k (\psi(\tau_{m1}) - \psi(\tau_{m1} + \tau_{m2})) \right) - z_{nk} \mathbb{E}_{\mathbf{v}} \left[\log \left(1 - \prod_{m=1}^k v_m \right) \right] \\ &\quad - \left(\frac{\rho_{ns}}{2\rho_{nr}} \right) \left(-2z_{nk} \bar{\phi}_k \mathbf{X}_{n\cdot}^T + z_{nk} (tr(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T) + 2z_{nk} \bar{\phi}_k \sum_{k' \neq k} \nu_{nk'} \bar{\phi}_{k'}^T \right)\end{aligned}\tag{37}$$

From the canonical parameterisation of Bernoulli distribution, we can get
(My understanding is that, assume a Bernoulli distribution $p(z; \nu) = \nu^z (1 - \nu)^{1-z}$, $\log(p) = z(\log \frac{\nu}{1-\nu})$)

$$\begin{aligned}\log \frac{\nu_{nk}}{1 - \nu_{nk}} &= \sum_{m=1}^k (\psi(\tau_{m1}) - \psi(\tau_{m1} + \tau_{m2})) - \mathbb{E}_{\mathbf{v}} \left[\log \left(1 - \prod_{m=1}^k v_m \right) \right] \\ &\quad - \left(\frac{\rho_{ns}}{2\rho_{nr}} \right) \left(-2\bar{\phi}_k \mathbf{X}_{n\cdot}^T + (tr(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T) + 2\bar{\phi}_k \sum_{k' \neq k} \nu_{nk'} \bar{\phi}_{k'}^T \right) \\ &\equiv \vartheta\end{aligned}\tag{38}$$

Gives the update

$$\nu_{nk} = \frac{1}{1 + e^{-\vartheta}}\tag{39}$$

3. Updates the feature probability τ

Identify terms that contain τ from the lower bound in equation 10

$$\begin{aligned}\mathcal{L}_\tau &= \sum_{k=1}^K [\log \alpha + (\alpha - 1)(\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2}))] \\ &\quad + \sum_{k=1}^K \sum_{n=1}^N \left[\nu_{nk} \left(\sum_{m=1}^k (\psi(\tau_{m1}) - \psi(\tau_{m1} + \tau_{m2})) \right) + (1 - \nu_{nk}) \mathbb{E}_{\mathbf{v}} [\log(1 - \prod_{m=1}^k v_m)] \right] \\ &\quad + \sum_{k=1}^K \left[\log \left(\frac{\Gamma(\tau_{k1}) \Gamma(\tau_{k2})}{\Gamma(\tau_{k1} + \tau_{k2})} \right) - (\tau_{k1} - 1) \psi(\tau_{k1}) - (\tau_{k2} - 1) \psi(\tau_{k2}) + (\tau_{k1} + \tau_{k2} - 2) \psi(\tau_{k1} + \tau_{k2}) \right]\end{aligned}\tag{40}$$

The expectation is bounded using the multinomial approximation

$$\begin{aligned}
\mathbb{E}_{\mathbf{v}}[\log(1 - \prod_{m=1}^k v_m)] \geq & \left(\sum_{m=1}^k q_{km} \psi(\tau_{m2}) \right) + \left(\sum_{m=1}^{k-1} \left(\sum_{j=m+1}^k q_{kj} \psi(\tau_{m1}) \right) \right) \\
& - \left(\sum_{m=1}^k \sum_{j=m}^k q_{kj} \psi(\tau_{m1} + \tau_{m2}) \right) - \sum_{m=1}^k q_{km} \log(q_{km})
\end{aligned} \tag{41}$$

Substitute into equation 40,

$$\begin{aligned}
\mathcal{L}_\tau = & \sum_{k=1}^K [\log \alpha + (\alpha - 1)(\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2}))] \\
& + \sum_{k=1}^K \sum_{n=1}^N \left[\nu_{nk} \left(\sum_{m=1}^k (\psi(\tau_{m1}) - \psi(\tau_{m1} + \tau_{m2})) \right) \right. \\
& + (1 - \nu_{nk}) \left(\left(\sum_{m=1}^k q_{km} \psi(\tau_{m2}) \right) + \left(\sum_{m=1}^{k-1} \left(\sum_{j=m+1}^k q_{kj} \psi(\tau_{m1}) \right) \right) \right. \\
& \left. \left. - \left(\sum_{m=1}^k \sum_{j=m}^k q_{kj} \psi(\tau_{m1} + \tau_{m2}) \right) - \sum_{m=1}^k q_{km} \log(q_{km}) \right) \right] \\
& + \sum_{k=1}^K \left[\log \left(\frac{\Gamma(\tau_{k1}) \Gamma(\tau_{k2})}{\Gamma(\tau_{k1} + \tau_{k2})} \right) - (\tau_{k1} - 1) \psi(\tau_{k1}) - (\tau_{k2} - 1) \psi(\tau_{k2}) + (\tau_{k1} + \tau_{k2} - 2) \psi(\tau_{k1} + \tau_{k2}) \right]
\end{aligned} \tag{42}$$

Remove terms that don't contain τ_k . We need to break down the summations. It's tricky especially for the second term. Since k is a constant now, it cannot be used as an iterative index. We change the notation first, replace k with m , replace m with i , then do the cleaning

$$\begin{aligned}
& \sum_{k=1}^K \sum_{n=1}^N \left[\nu_{nk} \left(\sum_{m=1}^k (\psi(\tau_{m1}) - \psi(\tau_{m1} + \tau_{m2})) \right) + (1 - \nu_{nk}) \left(\left(\sum_{m=1}^k q_{km} \psi(\tau_{m2}) \right) + \left(\sum_{m=1}^{k-1} \left(\sum_{j=m+1}^k q_{kj} \psi(\tau_{m1}) \right) \right) \right. \right. \\
& \quad \left. \left. - \left(\sum_{m=1}^k \sum_{j=m}^k q_{kj} \psi(\tau_{m1} + \tau_{m2}) \right) - \sum_{m=1}^k q_{km} \log(q_{km}) \right) \right] \\
& \xrightarrow{\text{renotate}} \sum_{m=1}^K \sum_{n=1}^N \left[\nu_{nm} \left(\sum_{i=1}^m (\psi(\tau_{i1}) - \psi(\tau_{i1} + \tau_{i2})) \right) + (1 - \nu_{nm}) \left(\left(\sum_{i=1}^m q_{mi} \psi(\tau_{i2}) \right) + \left(\sum_{i=1}^{m-1} \left(\sum_{j=i+1}^m q_{mj} \psi(\tau_{i1}) \right) \right) \right. \right. \\
& \quad \left. \left. - \left(\sum_{i=1}^m \sum_{j=i}^m q_{mj} \psi(\tau_{i1} + \tau_{i2}) \right) - \sum_{i=1}^m q_{mi} \log(q_{mi}) \right) \right]
\end{aligned} \tag{43}$$

To get all terms, the bottleneck is at the $\sum_{i=1}^{m-1}$, which means m should be $\geq k+1$ for it, and $\geq k$ for other terms

$$\begin{aligned} \xrightarrow{m \geq k \text{ or } k+1} & \sum_{m=k}^K \sum_{n=1}^N \left[\nu_{nm} \left(\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2}) \right) + (1 - \nu_{nm}) \left(q_{mk} \psi(\tau_{k2}) \right) \right. \\ & \left. - \left(\sum_{j=k}^m q_{mj} \psi(\tau_{k1} + \tau_{k2}) \right) - q_{mk} \log(q_{mk}) \right] + \sum_{m=k+1}^K \sum_{n=1}^N (1 - \nu_{nm}) \left(\sum_{j=k+1}^m q_{mj} \psi(\tau_{k1}) \right) \end{aligned} \quad (44)$$

Therefore,

$$\begin{aligned} \mathcal{L}_{\tau_k} = & \left[(\alpha - 1) + \sum_{m=k}^K \sum_{n=1}^N \nu_{nm} + \sum_{m=k+1}^K \sum_{n=1}^N (1 - \nu_{nm}) \left(\sum_{j=k+1}^m q_{mj} \right) - (\tau_{k1} - 1) \right] (\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2})) \\ & + \left[\sum_{m=k}^K \sum_{n=1}^N (1 - \nu_{nm}) q_{mk} - (\tau_{k2} - 1) \right] (\psi(\tau_{k2}) - \psi(\tau_{k1} + \tau_{k2})) + c \end{aligned} \quad (45)$$

gives the update formula

$$\begin{aligned} \tau_{k1} = & \alpha + \sum_{m=k}^K \sum_{n=1}^N \nu_{nm} + \sum_{m=k+1}^K \sum_{n=1}^N (1 - \nu_{nm}) \left(\sum_{j=k+1}^m q_{mj} \right) \\ \tau_{k2} = & \sum_{m=k}^K \sum_{n=1}^N (1 - \nu_{nm}) q_{mk} + 1 \end{aligned} \quad (46)$$

4. Updates the feature variance σ_A^2

$$\begin{aligned} \log q_{\rho_A}(\sigma_A^2) = & \mathbb{E}_{\mathbf{A}}[\log p(\mathbf{W}, \mathbf{X} | \boldsymbol{\theta})] + c \\ = & \mathbb{E}_{\mathbf{A}}[\log p(\sigma_A^2 | \boldsymbol{\sigma}_{A0}) + \sum_{k=1}^K \log p(\mathbf{A}_{k\cdot} | 0, \sigma_A^2 I)] + c \\ = & \log\left(\frac{\sigma_{Ar}^{\sigma_{As}}}{\Gamma(\sigma_{As})} \sigma_A^2 \sigma_{As}^{-1} e^{-\sigma_{Ar} \sigma_A^2}\right) + \sum_{k=1}^K \left(\frac{D}{2} \log\left(\frac{\sigma_A^2}{2\pi}\right) - \frac{\sigma_A^2}{2} \mathbb{E}_{\mathbf{A}}[\mathbf{A}_{k\cdot}^T \mathbf{A}_{k\cdot}] \right) + c \\ = & (\sigma_{As} - 1) \log(\sigma_A^2) - \sigma_{Ar} \sigma_A^2 + \frac{DK}{2} \log(\sigma_A^2) - \frac{\sigma_A^2}{2} \sum_{k=1}^K \left(\text{tr}(\boldsymbol{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T \right) + c \\ = & \left(\sigma_{As} + \frac{DK}{2} - 1 \right) \log(\sigma_A^2) - \left(\sigma_{Ar} + \frac{1}{2} \sum_{k=1}^K \left(\text{tr}(\boldsymbol{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T \right) \right) \sigma_A^2 + c \end{aligned} \quad (47)$$

Therefore, equations for updating σ_A^2 is

$$\begin{aligned}\rho_{As} &= \sigma_{As} + \frac{DK}{2} \\ \rho_{Ar} &= \sigma_{Ar} + \frac{1}{2} \sum_{k=1}^K \left(\text{tr}(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T \right)\end{aligned}\tag{48}$$

5-6. Similarly, σ_n^2 and α can be updated.

$$\begin{aligned}\log q_{\rho_n}(\sigma_n^2) &= \mathbb{E}_{\mathbf{Z}, \mathbf{A}}[\log p(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta})] + c \\ &= \mathbb{E}_{\mathbf{Z}, \mathbf{A}} \left[\log p(\sigma_n^2 | \sigma_{n0}) + \sum_{n=1}^N \log p(\mathbf{X}_{n\cdot} | \mathbf{Z}_{n\cdot}, \mathbf{A}, \sigma_n^2 I) \right] + c \\ &= \log \left(\frac{\sigma_{nr}^{\sigma_{ns}}}{\Gamma(\sigma_{ns})} \sigma_n^{2\sigma_{ns}-1} e^{-\sigma_{nr}\sigma_n^2} \right) + \sum_{n=1}^N \left(\frac{D}{2} \log \left(\frac{\sigma_n^2}{2\pi} \right) \right. \\ &\quad \left. - \frac{1}{2} \sigma_n^2 \left(\mathbf{X}_{n\cdot} \mathbf{X}_{n\cdot}^T - 2 \sum_{k=1}^K \nu_{nk} \bar{\phi}_k \mathbf{X}_{n\cdot}^T + 2 \sum_{k < k'} \nu_{nk} \nu_{nk'} \bar{\phi}_k \bar{\phi}_{k'}^T + \sum_{k=1}^K \nu_{nk} (\text{tr}(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T) \right) \right) \\ &= (\sigma_{ns} - 1) \log(\sigma_n^2) - \sigma_{nr} \sigma_n^2 + \frac{DN}{2} \log(\sigma_n^2) \\ &\quad - \frac{1}{2} \sigma_n^2 \sum_{n=1}^N \left(\mathbf{X}_{n\cdot} \mathbf{X}_{n\cdot}^T - 2 \sum_{k=1}^K \nu_{nk} \bar{\phi}_k \mathbf{X}_{n\cdot}^T + 2 \sum_{k < k'} \nu_{nk} \nu_{nk'} \bar{\phi}_k \bar{\phi}_{k'}^T + \sum_{k=1}^K \nu_{nk} (\text{tr}(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T) \right)\end{aligned}\tag{49}$$

Therefore, equations for updating σ_n^2 is

$$\begin{aligned}\rho_{ns} &= \sigma_{ns} + \frac{DN}{2} \\ \rho_{nr} &= \sigma_{nr} + \frac{1}{2} \sum_{n=1}^N \left(\mathbf{X}_{n\cdot} \mathbf{X}_{n\cdot}^T - 2 \sum_{k=1}^K \nu_{nk} \bar{\phi}_k \mathbf{X}_{n\cdot}^T + 2 \sum_{k < k'} \nu_{nk} \nu_{nk'} \bar{\phi}_k \bar{\phi}_{k'}^T + \sum_{k=1}^K \nu_{nk} (\text{tr}(\mathbf{\Phi}_k) + \bar{\phi}_k \bar{\phi}_k^T) \right)\end{aligned}\tag{50}$$

As to α

$$\begin{aligned}
\log q_{\rho_\alpha}(\alpha) &= \mathbb{E}_{\mathbf{v}}[\log p(\mathbf{W}, \mathbf{X}|\boldsymbol{\theta})] + c \\
&= \mathbb{E}_{\mathbf{v}}[\log p(\alpha|\alpha_0) + \sum_{k=1}^K \log p(v_k|\alpha)] + c \\
&= \log\left(\frac{\alpha_r^{\alpha_s}}{\Gamma(\alpha_s)} \alpha^{\alpha_s-1} e^{-\alpha_r \alpha}\right) + \sum_{k=1}^K \left(\log(\alpha) + (\alpha-1) \mathbb{E}_{\mathbf{v}} \log(v_k) \right) + c \\
&= (\alpha_s - 1) \log(\alpha) - \alpha_r \alpha + K \log(\alpha) + (\alpha - 1) \sum_{k=1}^K (\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2})) + c
\end{aligned} \tag{51}$$

Equations for updating α is

$$\begin{aligned}
\rho_{\alpha s} &= \alpha_s + K \\
\rho_{\alpha r} &= \alpha_r - (\alpha - 1) \sum_{k=1}^K (\psi(\tau_{k1}) - \psi(\tau_{k1} + \tau_{k2}))
\end{aligned} \tag{52}$$

Only concern is that $\rho_{\alpha r}$ may go below 0

References

- [1] Matthew James Beal. *Variational algorithms for approximate Bayesian inference*. University of London, University College London (United Kingdom), 2003.
- [2] David M Blei and Michael I Jordan. “Variational methods for the Dirichlet process”. In: *Proceedings of the twenty-first international conference on Machine learning*. 2004, p. 12.
- [3] Finale Doshi et al. “Variational inference for the Indian buffet process”. In: *Artificial Intelligence and Statistics*. PMLR. 2009, pp. 137–144.
- [4] Diederik P Kingma and Max Welling. “Auto-encoding variational bayes”. In: *arXiv preprint arXiv:1312.6114* (2013).
- [5] Stephen Odaibo. “Tutorial: Deriving the standard variational autoencoder (VAE) loss function”. In: *arXiv preprint arXiv:1907.08956* (2019).