

I will start off with some useful properties I need to use later on in the ELBO.

## 1 Negative cross entropies

### 1.1 Two Normals

Note: All the normals are parametrized using the precision matrix.

$$q \sim \mathcal{N}(x|m, L)$$

$$p \sim \mathcal{N}(x|\mu, \Lambda)$$

$$\begin{aligned} \int q(x) \ln p(x) dx &= \int \mathcal{N}(x|m, L) \left( -\frac{K}{2} \ln 2\pi + \frac{1}{2} \ln |\Lambda| - \frac{1}{2} \left( \text{Tr} \Lambda \{ (x - \mu)(x - \mu)^T \} \right) \right) dx \\ &= -\frac{K}{2} \ln 2\pi + \frac{1}{2} \ln |\Lambda| + \int \mathcal{N}(x|m, L) \left( -\frac{1}{2} \left( \text{Tr} \Lambda \{ (x - \mu)(x - \mu)^T \} \right) \right) dx \\ &= -\frac{K}{2} \ln 2\pi + \frac{1}{2} \ln |\Lambda| + \int \mathcal{N}(x|m, L) \left( -\frac{1}{2} \left( \text{Tr} \Lambda \{ xx^T + \mu\mu^T - x\mu^T - \mu x^T \} \right) \right) dx \end{aligned}$$

$$\text{We should note that } \mathbb{E}_q [xx^T] = \text{Cov}_q + \mathbb{E}_q [x] \mathbb{E}_q [x]^T$$

$$\mathbb{E}_q [x] = m \text{ and } \text{Cov}_q = L^{-1}$$

$$\begin{aligned} \int \mathcal{N}(x|m, L) \left( -\frac{1}{2} \left( \text{Tr} \left[ \Lambda \{ xx^T + \mu\mu^T - x\mu^T - \mu x^T \} \right] \right) \right) dx &= -\frac{1}{2} \text{Tr} \left[ (\Lambda L^{-1} + \Lambda m m^T) + \Lambda (m m^T - \mu m^T - m \mu^T) \right] \\ &= -\frac{1}{2} \left( \text{Tr} [\Lambda L^{-1}] + (m - \mu)^T \Lambda (m - \mu) \right) \end{aligned}$$

Hence we have:

$$\boxed{\mathbb{E}_q [\ln p(x)] = -\frac{K}{2} \ln 2\pi + \frac{1}{2} \ln |\Lambda| - \frac{1}{2} \left( \text{Tr} [\Lambda L^{-1}] + (m - \mu)^T \Lambda (m - \mu) \right)}$$

### 1.2 Two Wisharts

$$\Lambda \sim q \sim \mathcal{W}(v, W)$$

$$\Lambda \sim p \sim \mathcal{W}(n, S)$$

$$\begin{aligned} \int q(\Lambda) \ln p(\Lambda) d\Lambda &= \mathbb{E}_q [\ln p(\Lambda)] \\ &= \mathbb{E}_q \left[ \ln \frac{|\Lambda|^{\frac{n-K-1}{2}} \exp(-\frac{1}{2} \text{Tr}(S^{-1} \Lambda))}{2^{\frac{nK}{2}} |S|^{n/2} \Gamma_p(\frac{n}{2})} \right] \\ &= \mathbb{E}_q \left[ -\frac{nk}{2} \ln 2 - \frac{n}{2} \ln |S| - \ln \Gamma_K(\frac{n}{2}) \right. \\ &\quad \left. + \frac{n-K-1}{2} \ln |\Lambda| - \frac{1}{2} \text{Tr}(S^{-1} \Lambda) \right] \\ &= -\frac{nk}{2} \ln 2 - \frac{n}{2} \ln |S| - \ln \Gamma_K(\frac{n}{2}) \\ &\quad + \frac{n-K-1}{2} \left( \psi_K(\frac{v}{2}) + K \ln 2 + \ln |W| \right) - \frac{v}{2} \text{Tr}(S^{-1} W) \end{aligned}$$

Note that:

$$\mathbb{E}_q[\Lambda] = vW$$

$$\mathbb{E}_q[\ln |\Lambda|] = \psi_K(\frac{v}{2}) + K \ln 2 + \ln |W|$$

$$\psi_K(\frac{v}{2}) = \sum_{i:1}^K \psi(\frac{v-i+1}{2})$$

$$\ln \Gamma_K(\frac{n}{2}) = \frac{K(K-1)}{4} \ln \pi + \sum_{i:1}^K \ln \Gamma(\frac{n-i+1}{2})$$

$$\begin{aligned} \mathbb{E}_q[\ln p(\Lambda)] &= -\frac{K(K+1)}{2} \ln 2 + \frac{n-K-1}{2} \psi_K(\frac{v}{2}) - \ln \Gamma_K(\frac{n}{2}) \\ &\quad - \frac{v}{2} \text{Tr}(S^{-1}W) + \frac{n-K-1}{2} \ln |W| - \frac{n}{2} \ln |S| \end{aligned}$$

so we have:

$$\mathbb{E}_q[\ln p(\Lambda)] = -\frac{K(K+1)}{2} \ln 2 + \frac{n-K-1}{2} \psi_K(\frac{v}{2}) - \ln \Gamma_K(\frac{n}{2}) - \frac{v}{2} \text{Tr}(S^{-1}W) + \frac{n-K-1}{2} \ln |W| - \frac{n}{2} \ln |S|$$

or

$$\mathbb{E}_q[\ln p(\Lambda)] = -\frac{K(K+1)}{2} \ln 2 + \frac{n-K-1}{2} \psi_K(\frac{v}{2}) - \ln \Gamma_K(\frac{n}{2}) - \frac{v}{2} \text{Tr}(S^{-1}W) - \frac{K+1}{2} \ln |W| + \frac{n}{2} \ln |S^{-1}W|$$

### 1.3 Two Betas

$$\beta \sim q \sim \text{Beta}(b)$$

$$\beta \sim p \sim \text{Beta}(\eta)$$

$$\begin{aligned} \mathbb{E}_q[\ln p(\beta)] &= \mathbb{E}_q \left[ \ln \Gamma(\eta_0 + \eta_1) - \ln \Gamma(\eta_0) - \ln \Gamma(\eta_1) + (\eta_0 - 1) \ln \beta + (\eta_1 - 1) \ln (1 - \beta) \right] \\ &= \ln \Gamma(\eta_0 + \eta_1) - \ln \Gamma(\eta_0) - \ln \Gamma(\eta_1) + (\eta_0 - 1) (\psi(b_0) - \psi(b_0 + b_1)) + (\eta_1 - 1) (\psi(b_1) - \psi(b_0 + b_1)) \\ &= \ln \Gamma(\eta_0 + \eta_1) - \ln \Gamma(\eta_0) - \ln \Gamma(\eta_1) + (\eta_0 - 1) \psi(b_0) + (\eta_1 - 1) \psi(b_1) - (\eta_0 + \eta_1 - 2) \psi(b_0 + b_1) \end{aligned}$$

$$\text{Note that } \mathbb{E}_q[\ln \beta] = \psi(b_0) - \psi(b_0 + b_1)$$

so :

$$\mathbb{E}_q[\ln p(\beta)] = \ln \Gamma(\eta_0 + \eta_1) - \ln \Gamma(\eta_0) - \ln \Gamma(\eta_1) + (\eta_0 - 1) \psi(b_0) + (\eta_1 - 1) \psi(b_1) - (\eta_0 + \eta_1 - 2) \psi(b_0 + b_1)$$

## 2 Entropies

### 2.1 Normal

$$q(x) \sim \mathcal{N}(m, M)$$

$$H[q] = \frac{K}{2} \ln(2\pi) + \frac{K}{2} - \frac{1}{2} \ln |M|$$

### 2.2 Wishart

$$\Lambda \sim q \sim \mathcal{W}(v, W)$$

$$\begin{aligned} H[q] &= -\frac{v-K-1}{2} \mathbb{E}_q \ln |\Lambda| - (-\frac{1}{2} \mathbb{E}_q \text{Tr}(W^{-1}\Lambda)) + \frac{v}{2} \ln |W| + \frac{vK}{2} \ln 2 + \ln \Gamma_K(\frac{v}{2}) \\ &= -\frac{v-K-1}{2} (\psi_K(\frac{v}{2}) + \frac{Kv}{2} + K \ln 2 + \ln |W|) + \frac{v}{2} \ln |W| + \frac{vK}{2} \ln 2 + \ln \Gamma_K(\frac{v}{2}) \\ &= \frac{K(K+1)}{2} \ln 2 + \frac{K+1}{2} \ln |W| - \frac{v-K-1}{2} \psi_K(\frac{v}{2}) + \ln \Gamma_K(\frac{v}{2}) + \frac{Kv}{2} \end{aligned}$$

so

$$H[q] = \frac{K(K+1)}{2} \ln 2 + \frac{K+1}{2} \ln |W| - \frac{v-K-1}{2} \psi_K(\frac{v}{2}) + \ln \Gamma_K(\frac{v}{2}) + \frac{Kv}{2}$$

### 2.3 Beta

$$\beta \sim q \sim \text{Beta}(b)$$

$$\begin{aligned} H[q] &= \ln \Gamma(b_0) + \ln \Gamma(b_1) - \ln \Gamma(b_0 + b_1) - (b_0 - 1) \mathbb{E}_q[\ln \beta] - (b_1 - 1) \mathbb{E}_q[\ln (1 - \beta)] \\ &= \ln \Gamma(b_0) + \ln \Gamma(b_1) - \ln \Gamma(b_0 + b_1) - (b_0 - 1) \psi(b_0) - (b_1 - 1) \psi(b_1) + (b_0 + b_1 - 2) \psi(b_0 + b_1) \end{aligned}$$

So,

$$H[q] = \ln \Gamma(b_0) + \ln \Gamma(b_1) - \ln \Gamma(b_0 + b_1) - (b_0 - 1)\psi(b_0) - (b_1 - 1)\psi(b_1) + (b_0 + b_1 - 2)\psi(b_0 + b_1)$$

## 2.4 Multinomial(,1) or Categorical

$$z \sim q \sim \text{Cat}(\phi)$$

$$H[q] = - \sum_k \mathbb{E}_q[z_k] \ln \phi_k$$

so,

$$H[q] = - \sum_k \phi_k \ln \phi_k$$

## 3 Variational ELBO

$$\mathcal{L} = \mathbb{E}_q[\ln p(\text{joint})] + H_q[\text{params}]$$

$$\begin{aligned} \ln p(\text{joint}) &= \ln p(\mu|m_0, M_0) + \ln p(\Lambda|\ell_0, L_0) + \sum_a \ln p(\theta_a|\mu, \Lambda) + \sum_a \sum_b \ln p(z_{a \rightarrow b}|\theta_a) \\ &\quad + \sum_a \sum_b \ln p(z_{a \leftarrow b}|\theta_b) + \sum_k \ln p(\beta_{kk}|\eta) + \sum_a \sum_b \ln p(y_{ab}|z_{a \rightarrow b}, z_{a \leftarrow b}, \beta) \end{aligned}$$

$$H_q[\text{params}] = H_q[\mu] + H_q[\Lambda] + H_q[\theta] + H_q[\beta] + H_q[z_{\rightarrow}] + H_q[z_{\leftarrow}]$$

Furthermore,

$$\begin{aligned} \mathbb{E}_q[\ln p(\text{joint})] &= \mathbb{E}_q[\ln p(\mu|m_0, M_0)] + \mathbb{E}_q[\ln p(\Lambda|\ell_0, L_0)] + \sum_a \mathbb{E}_q[\ln p(\theta_a|\mu, \Lambda)] + \sum_a \sum_b \mathbb{E}_q[\ln p(z_{a \rightarrow b}|\theta_a)] \\ &\quad + \sum_a \sum_b \mathbb{E}_q[\ln p(z_{a \leftarrow b}|\theta_b)] + \sum_k \mathbb{E}_q[\ln p(\beta_{kk}|\eta)] + \sum_a \sum_b \mathbb{E}_q[\ln p(y_{ab}|z_{a \rightarrow b}, z_{a \leftarrow b}, \beta)] \end{aligned}$$

We parametrize the variational distribution as follows:

$$\begin{aligned} \mu &\sim q(\mu|m, M) \sim \mathcal{N}(\mu|m, M) \\ \Lambda &\sim q(\Lambda|\ell, L) \sim \mathcal{W}(\Lambda|\ell, L) \\ \theta_a &\sim q(\theta_a|\mu_a, \Lambda_a) \sim \mathcal{N}(\theta_a|\mu_a, \Lambda_a) \\ \beta_{kk} &\sim q(\beta_{kk}|b_k) \sim \mathcal{B}(b_{k0}, b_{k1}) \\ z_{a \rightarrow b} &\sim q(z_{a \rightarrow b}|\phi_{a \rightarrow b}) \sim \text{Cat}(z_{a \rightarrow b}|\phi_{a \rightarrow b}) \\ z_{a \leftarrow b} &\sim q(z_{a \leftarrow b}|\phi_{a \leftarrow b}) \sim \text{Cat}(z_{a \leftarrow b}|\phi_{a \leftarrow b}) \end{aligned}$$

Using the results from above regarding the negative cross entropies:

$$\begin{aligned}
\mathbb{E}_q[\ln p(\text{joint})] &= -\frac{K}{2} \ln 2\pi + \frac{1}{2} \ln |M_0| - \frac{1}{2} \left( \text{Tr} M_0 \left[ M^{-1} + (m - m_0)(m - m_0)^T \right] \right) \\
&\quad - \frac{K(K+1)}{2} \ln 2 + \frac{\ell_0 - K - 1}{2} \psi_K\left(\frac{\ell}{2}\right) - \ln \Gamma_K\left(\frac{\ell_0}{2}\right) - \frac{\ell}{2} \text{Tr} (L_0^{-1} L) - \frac{K+1}{2} \ln |L| + \frac{\ell_0}{2} \ln |L_0^{-1} L| \\
&\quad - \sum_a \frac{K}{2} \ln 2\pi + \frac{1}{2} \sum_a \psi_K\left(\frac{\ell}{2}\right) + \frac{1}{2} \sum_a K \ln 2 + \frac{1}{2} \sum_a \ln |L| \\
&\quad - \frac{\ell}{2} \left( \text{Tr} \left[ L \left( \sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) + \sum_a M^{-1} \right) \right] \right) \\
&\quad + \sum_a \sum_b \sum_k \phi_{a \rightarrow b, k} \mu_{a, k} - \sum_a \sum_b \mathbb{E}_q[\ln (\sum_l \exp(\theta_{a, l}))] \\
&\quad + \sum_a \sum_b \sum_k \phi_{a \leftarrow b, k} \mu_{b, k} - \sum_a \sum_b \mathbb{E}_q[\ln (\sum_l \exp(\theta_{b, l}))] \\
&\quad + \sum_k \ln \Gamma(\eta_0 + \eta_1) - \sum_k \ln \Gamma(\eta_0) - \sum_k \ln \Gamma(\eta_1) + \sum_k (\eta_0 - 1) \psi(b_{k0}) \\
&\quad + \sum_k (\eta_1 - 1) \psi(b_{k1}) - \sum_k (\eta_0 + \eta_1 - 2) \psi(b_{k0} + b_{k1}) \\
&\quad + \sum_{a, b \in \text{link}} \sum_k \phi_{a \rightarrow b, k} \phi_{a \leftarrow b, k} (\psi(b_{k0}) - \psi(b_{k0} + b_{k1}) - \ln \epsilon) + \ln \epsilon \\
&\quad + \sum_{a, b \notin \text{link}} \sum_k \phi_{a \rightarrow b, k} \phi_{a \leftarrow b, k} (\psi(b_{k1}) - \psi(b_{k0} + b_{k1}) - \ln(1 - \epsilon)) + \ln(1 - \epsilon)
\end{aligned}$$

$$\mathbb{E}_q[\Lambda] = \ell L$$

$$\mathbb{E}_q[\ln |\Lambda|] = \psi_K\left(\frac{\ell}{2}\right) + K \ln 2 + \ln |L|$$

$$\begin{aligned}
& - \sum_a \frac{K}{2} \ln 2\pi + \sum_a \frac{1}{2} \mathbb{E}_q \left\{ \ln |\Lambda| \right\} \\
& - \sum_a \frac{1}{2} \left( \text{Tr} \left[ \mathbb{E}_q \left\{ \Lambda \right\} \Lambda_a^{-1} \right] + \mathbb{E}_q \left\{ (\mu_a - \mu)^T \Lambda (\mu_a - \mu) \right\} \right) = \\
& \quad - \sum_a \frac{K}{2} \ln 2\pi + \sum_a \psi_K\left(\frac{\ell}{2}\right) + \sum_a K \ln 2 + \sum_a \ln |L| \\
& \quad - \frac{\ell}{2} \left( \text{Tr} \left[ L \left( \sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) \right) \right] \right)
\end{aligned}$$

For the expression  $\mathbb{E}_q[\ln (\sum_l \exp(\theta_{a, l}))]$ , we use the Jensen's inequality to acquire:

$$\begin{aligned}
\mathbb{E}_q[\ln (\sum_l \exp(\theta_{a, l}))] &\leq \ln (\sum_l \mathbb{E}_q[\exp(\theta_{a, l})]) \\
&= \ln (\sum_l \exp(\mu_{a, l} + \frac{1}{2} \text{diag}(\Lambda_a^{-1})_{ll}))
\end{aligned}$$

We can introduce another bound that introduces a new variational parameter per individual:

$$\mathbb{E}_q[\ln (\sum_l \exp(\theta_{a, l}))] \leq \zeta_a^{-1} \sum_l \exp(\mu_{a, l} + \frac{1}{2} \text{diag}(\Lambda_a^{-1})_{ll}) + \ln \zeta_a - 1$$

Moreover, using the entropies from above:

$$\begin{aligned}
H_q[params] = & \frac{K}{2} \ln(2\pi) + \frac{K}{2} - \frac{1}{2} \ln|M| \\
& + \frac{K(K+1)}{2} \ln 2 + \frac{K+1}{2} \ln|L| - \frac{\ell-K-1}{2} \psi_K\left(\frac{\ell}{2}\right) + \ln \Gamma_K\left(\frac{\ell}{2}\right) + \frac{K\ell}{2} \\
& + \sum_a \frac{K}{2} \ln(2\pi) + \sum_a \frac{K}{2} - \sum_a \frac{1}{2} \ln|\Lambda_a| \\
& + \sum_k \ln \Gamma(b_{k0}) + \sum_k \ln \Gamma(b_{k1}) - \sum_k \ln \Gamma(b_{k0} + b_{k1}) - \sum_k (b_{k0} - 1) \psi(b_{k0}) \\
& - \sum_k (b_{k1} - 1) \psi(b_{k1}) + \sum_k (b_{k0} + b_{k1} - 2) \psi(b_{k0} + b_{k1}) \\
& - \sum_a \sum_b \sum_k \phi_{a \rightarrow b, k} \ln \phi_{a \rightarrow b, k} \\
& - \sum_a \sum_b \sum_k \phi_{a \leftarrow b, k} \ln \phi_{a \leftarrow b, k}
\end{aligned}$$

Note that here I assume the following for the hyperparameters:

$$\begin{aligned}
m_0 &= 0 \\
M_0 &= I \\
\ell_0 &= K \\
L_0 &= \frac{1}{K} I \\
\eta_0 &> 1 \\
\eta_1 &= 1
\end{aligned}$$

Finally, we have the following:

$$\begin{aligned}
\mathcal{L} = & -\frac{1}{2} \left( K \ln 2\pi + \text{tr}(mm^T) + \text{tr} M^{-1} \right) \\
& -\frac{1}{2} \left( -K^2 \ln K + \ln |L| + \ell K + \text{tr} L + \frac{K(K-1)}{2} \ln \pi + \right. \\
& 2 \sum_i \ln \Gamma\left(\frac{K-i+1}{2}\right) + \sum_i \Psi\left(\frac{\ell-i+1}{2}\right) + K(K+1) \ln 2 \Big) \\
& -\frac{1}{2} \left( K \ln 2\pi - \sum_i \Psi\left(\frac{\ell-i+1}{2}\right) - K \ln 2 - \ln |L| + \right. \\
& \ell \text{tr} \left\{ L[(\mu_a - m)(\mu_a - m)^T + M^{-1} + \Lambda_a^{-1}] \right\} \Big) \\
& + \sum_a \sum_{b \in \text{sink}(a)} \left( \sum_k \phi_{a \rightarrow b, k} \mu_{a, k} - \ln \sum_l \exp(\mu_{a, l} + \frac{1}{2} \Lambda_{a, l}^{-1}) \right) \\
& + \sum_a \sum_{b \notin \text{sink}(a)} \left( \sum_k \phi_{a \rightarrow b, k} \mu_{a, k} - \ln \sum_l \exp(\mu_{a, l} + \frac{1}{2} \Lambda_{a, l}^{-1}) \right) \\
& + \sum_a \sum_{b \in \text{source}(a)} \left( \sum_k \phi_{b \leftarrow a, k} \mu_{a, k} - \ln \sum_l \exp(\mu_{a, l} + \frac{1}{2} \Lambda_{a, l}^{-1}) \right) \\
& + \sum_a \sum_{b \notin \text{source}(a)} \left( \sum_k \phi_{b \leftarrow a, k} \mu_{a, k} - \ln \sum_l \exp(\mu_{a, l} + \frac{1}{2} \Lambda_{a, l}^{-1}) \right) \\
& + \sum_k (\eta_0 - 1) \Psi(b_{k0}) - \sum_k (\eta_0 - 2) \Psi(b_{k0} + b_{k1}) \\
& + \sum_a \sum_{b \in \text{sink}(a)} \sum_k \left( \phi_{a \rightarrow b, k} \phi_{a \leftarrow b, k} (\Psi(b_{k0}) - \Psi(b_{k0} + b_{k1}) - \ln \epsilon) + \ln \epsilon \right) \\
& + \sum_a \sum_{b \notin \text{sink}(a)} \sum_k \left( \phi_{a \rightarrow b, k} \phi_{a \leftarrow b, k} (\Psi(b_{k1}) - \Psi(b_{k0} + b_{k1}) - \ln(1 - \epsilon)) + \ln(1 - \epsilon) \right) \\
& + \frac{1}{2} \left( K \ln 2\pi + K - \ln |M| \right) \\
& + \frac{1}{2} \left( (K+1) \ln |L| + K(K+1) \ln 2 + \ell K + \frac{1}{2} K(K-1) \ln \pi \right. \\
& + 2 \sum_i \ln \Gamma\left(\frac{\ell-i+1}{2}\right) - \frac{\ell-K-1}{2} \sum_i \Psi\left(\frac{\ell-i+1}{2}\right) \Big) \\
& + \frac{1}{2} \sum_a \left( K \ln 2\pi - \ln |\Lambda_a| + K \right) \\
& + \sum_k \left( \ln \Gamma(b_{k0}) + \ln \Gamma(b_{k1}) - \ln \Gamma(b_{k0} + b_{k1}) - (b_{k0} - 1) \Psi(b_{k0}) - \right. \\
& (b_{k1} - 1) \Psi(b_{k1}) + (b_{k0} + b_{k1} - 2) \Psi(b_{k0} + b_{k1}) \Big) \\
& - \sum_a \sum_{b \in \text{sink}(a)} \sum_k \left( \phi_{a \rightarrow b, k} \ln \phi_{a \rightarrow b, k} \right) \\
& - \sum_a \sum_{b \notin \text{sink}(a)} \sum_k \left( \phi_{a \rightarrow b, k} \ln \phi_{a \rightarrow b, k} \right) \\
& - \sum_a \sum_{b \in \text{sink}(a)} \sum_k \left( \phi_{a \leftarrow b, k} \ln \phi_{a \leftarrow b, k} \right) \\
& - \sum_a \sum_{b \notin \text{sink}(a)} \sum_k \left( \phi_{a \leftarrow b, k} \ln \phi_{a \leftarrow b, k} \right)
\end{aligned}$$

## 4 ELBO Gradients

### 4.1 Gradient with respect to $m$

$$\begin{aligned}
\mathcal{L}_m &= -\frac{1}{2} \left( \text{Tr } m m^T \right) \\
&\quad - \frac{\ell}{2} \left( \text{Tr } L \left( \sum_a (m - \mu_a)(m - \mu_a)^T \right) \right) \\
&\propto \text{Tr } m m^T \\
&\quad + \ell \left( \text{Tr } L \left( \sum_a m m^T + \mu_a \mu_a^T - m \mu_a^T - \mu_a m^T \right) \right) \\
&= m^T \left( I + N \ell L \right) m - m^T \left( \ell L \sum_a \mu_a \right) - \left( \ell \sum_a \mu_a^T L \right) m \\
&\Rightarrow \\
\nabla_m \mathcal{L}_m &\propto (I + N \ell L) m - (\ell L \sum_a \mu_a) = 0 \\
&\Rightarrow \\
&\boxed{m = (I + N \ell L)^{-1} (\ell L \sum_a \mu_a)}
\end{aligned}$$

In minibatch node sampling this would be

$$\boxed{m = (I + N \ell L)^{-1} (\ell L \frac{N}{\#mbnodes} \sum_{a \in mbnodes} \mu_a)}$$

### 4.2 Gradient with respect to $M$

$$\begin{aligned}
\mathcal{L}_M &= -\frac{1}{2} \left( \text{Tr } M^{-1} \right) \\
&\quad - \frac{\ell}{2} \text{Tr } N L M^{-1} \\
&\quad - \frac{1}{2} \ln |M| \\
&\propto \text{Tr } M^{-1} + \ell \text{Tr } N L M^{-1} + \ln |M| \\
&\Rightarrow \\
\nabla_M \mathcal{L}_M &= -(M^{-1} M^{-1})^T - \ell N (M^{-1} L M^{-1})^T + (M^{-1})^T = 0 \\
&\quad \text{transpose} \\
&\quad M \times () \times M \\
&\quad \Leftrightarrow -I - N \ell L + M = 0 \\
&\quad \boxed{M = I + N \ell L}
\end{aligned}$$

### 4.3 Gradient with respect to $L$

$$\begin{aligned}
\mathcal{L}_L &= -\frac{\ell}{2} \text{Tr}(KIL) - \frac{K+1}{2} \ln |L| + \frac{K}{2} \ln |L_0^{-1}L| \\
&\quad + \frac{1}{2} \sum_a \ln |L| - \frac{\ell}{2} \left( \text{Tr} \left[ L \left( \sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) + \sum_a M^{-1} \right) \right] \right) \\
&\quad + \frac{K+1}{2} \ln |L| \\
&\propto -\ell \text{Tr}(KIL) - (K+1) \ln |L| + K \ln |KIL| \\
&\quad + \sum_a \ln |L| - \ell \left( \text{Tr} \left[ L \left( \sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) + \sum_a M^{-1} \right) \right] \right) \\
&\quad + (K+1) \ln |L| \\
&\Rightarrow \\
\nabla_L \mathcal{L}_L &= -\ell(KI)^T + K(L^{-1})^T + N(L^{-1})^T - \ell \left( \sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) + \sum_a M^{-1} \right)^T = 0 \\
&\quad \ell(KI + \sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) + \sum_a M^{-1}) = (N + K)L^{-1} \\
&\Rightarrow \boxed{L = \frac{N + K}{\ell} \left( (KI + \sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) + \sum_a M^{-1}) \right)^{-1}}
\end{aligned}$$

optimizing simultaeneously with  $\ell$  in the minibatch setting:

$$\boxed{L = \left( (KI + \frac{N}{\#mbnodes} \sum_{a \in mbnodes} (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T)) + N \times M^{-1} \right)^{-1}}$$



#### 4.4 Gradient with respect to $\ell$

$$\begin{aligned}
\mathcal{L}_\ell &= -\frac{1}{2}\psi_K\left(\frac{\ell}{2}\right) + \frac{1}{2}\sum_a \psi_K\left(\frac{\ell}{2}\right) - \frac{\ell}{2}\text{Tr}(KIL) \\
&\quad - \frac{\ell}{2}\left(\text{Tr}\left[L\left(\sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) + \sum_a M^{-1}\right)\right]\right) \\
&\quad - \frac{\ell-K-1}{2}\psi_K\left(\frac{\ell}{2}\right) + \ln \Gamma_K\left(\frac{\ell}{2}\right) + \frac{K\ell}{2} \\
&\propto (-1)\psi_K\left(\frac{\ell}{2}\right) + \sum_a \psi_K\left(\frac{\ell}{2}\right) - \ell\text{Tr}(KIL) \\
&\quad - \ell\left(\text{Tr}\left[L\left(\sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) + \sum_a M^{-1}\right)\right]\right) \\
&\quad - (\ell - K - 1)\psi_K\left(\frac{\ell}{2}\right) + 2\ln \Gamma_K\left(\frac{\ell}{2}\right) + K\ell \\
&\Rightarrow \left(\sum_{i:1}^K \psi\left(\frac{\ell-i+1}{2}\right)\right)\left(K - \mathbb{K} - \mathbb{I} + N - l + \mathbb{K} + \mathbb{I}\right) + 2\left(\frac{K(K-1)}{4}\ln \pi + \sum_{i:1}^K \ln \Gamma\left(\frac{\ell-i+1}{2}\right)\right) \\
&\quad + \ell(K - \text{Tr}\left[L\left(KI + \sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) + \sum_a M^{-1}\right)\right]) \\
&\propto \left(\sum_{i:1}^K \psi\left(\frac{\ell-i+1}{2}\right)\right)\left(K + N - \ell\right) + 2\left(\sum_{i:1}^K \ln \Gamma\left(\frac{\ell-i+1}{2}\right)\right) \\
&\quad + \ell(K - \text{Tr}\left[L\left(KI + \sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) + \sum_a M^{-1}\right)\right]) \\
&\Rightarrow \\
\nabla_\ell \mathcal{L}_\ell &= \frac{1}{2}\left(\sum_{i:1}^K \psi'\left(\frac{\ell-i+1}{2}\right)\right)\left(K + N - \ell\right) \\
&\quad + K - \text{Tr}\left[L\left(KI + \sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) + \sum_a M^{-1}\right)\right] = 0 \\
&\quad \text{simultaneously with } L \\
&\Rightarrow L\left(KI + \sum_a (\Lambda_a^{-1} + (m - \mu_a)(m - \mu_a)^T) + \sum_a M^{-1}\right) = I_K \\
&\text{hence, } K - \text{Tr } I_K = 0 \\
&\Rightarrow \boxed{\ell = K + N}
\end{aligned}$$

#### 4.5 Gradient with respect to $b_k$

$$\begin{aligned}
\mathcal{L}_{b_k} &= (\eta_0 - 1)\psi(b_{k0}) + -(\eta_0 - 2)\psi(b_{k0} + b_{k1}) \\
&+ \sum_{a,b \in link} \phi_{a \rightarrow b,k} \phi_{a \leftarrow b,k} (\psi(b_{k0}) - \psi(b_{k0} + b_{k1})) \\
&+ \sum_{a,b \notin link} \phi_{a \rightarrow b,k} \phi_{a \leftarrow b,k} (\psi(b_{k1}) - \psi(b_{k0} + b_{k1})) \\
&+ \ln \Gamma(b_{k0}) + \ln \Gamma(b_{k1}) - \ln \Gamma(b_{k0} + b_{k1}) - (b_{k0} - 1)\psi(b_{k0}) \\
&- (b_{k1} - 1)\psi(b_{k1}) + (b_{k0} + b_{k1} - 2)\psi(b_{k0} + b_{k1}) \\
&\text{simultaneously optimizing } b_{k0}, b_{k1} \\
\Rightarrow &\text{Similar to our previous results} \\
\nabla_{b_{k0}} \mathcal{L}_{b_k} &= 0 \\
\Rightarrow &\boxed{b_{k0} = \eta_0 + \frac{\#trainlinks}{\#mblinks} \sum_{a,b \in mblinks} \phi_{a \rightarrow b,k} \phi_{a \leftarrow b,k}} \\
\nabla_{b_{k1}} \mathcal{L}_{b_k} &= 0 \\
&\boxed{b_{k1} = 1 + \frac{\#trainnonlinks}{\#mbnonlinks} \sum_{a,b \notin mblinks} \phi_{a \rightarrow b,k} \phi_{a \leftarrow b,k}}
\end{aligned}$$

#### 4.6 Gradient with respect to $\phi_{a \rightarrow b,k}$ for links

$$\begin{aligned}
\mathcal{L}_{\phi_{a \rightarrow b,k}} &= \phi_{a \rightarrow b,k} \mu_{a,k} \\
&+ \phi_{a \rightarrow b,k} \phi_{a \leftarrow b,k} (\psi(b_{k0}) - \psi(b_{k0} + b_{k1}) - \ln \epsilon) \\
&- \phi_{a \rightarrow b,k} \ln \phi_{a \rightarrow b,k} \\
&= \phi_{a \rightarrow b,k} (\mu_{a,k} + \phi_{a \leftarrow b,k} (\psi(b_{k0}) - \psi(b_{k0} + b_{k1}) - \ln \epsilon) - \ln \phi_{a \rightarrow b,k}) \\
\nabla_{\phi_{a \rightarrow b,k}} \mathcal{L}_{\phi_{a \rightarrow b,k}} &= \mu_{a,k} + \phi_{a \leftarrow b,k} (\psi(b_{k0}) - \psi(b_{k0} + b_{k1}) - \ln \epsilon) - \ln \phi_{a \rightarrow b,k} = 0 \\
&\boxed{\phi_{a \rightarrow b,k} \propto \exp \left\{ \mu_{a,k} + \phi_{a \leftarrow b,k} (\psi(b_{k0}) - \psi(b_{k0} + b_{k1}) - \ln \epsilon) \right\}}
\end{aligned}$$

#### 4.7 Gradient with respect to $\phi_{a \leftarrow b,k}$ for links

$$\begin{aligned}
\mathcal{L}_{\phi_{a \leftarrow b,k}} &= \phi_{a \leftarrow b,k} \mu_{b,k} \\
&+ \phi_{a \rightarrow b,k} \phi_{a \leftarrow b,k} (\psi(b_{k0}) - \psi(b_{k0} + b_{k1}) - \ln \epsilon) \\
&- \phi_{a \leftarrow b,k} \ln \phi_{a \leftarrow b,k} \\
&= \phi_{a \leftarrow b,k} (\mu_{b,k} + \phi_{a \rightarrow b,k} (\psi(b_{k0}) - \psi(b_{k0} + b_{k1}) - \ln \epsilon) - \ln \phi_{a \leftarrow b,k}) \\
\nabla_{\phi_{a \rightarrow b,k}} \mathcal{L}_{\phi_{a \leftarrow b,k}} &= \mu_{b,k} + \phi_{a \rightarrow b,k} (\psi(b_{k0}) - \psi(b_{k0} + b_{k1}) - \ln \epsilon) - \ln \phi_{a \leftarrow b,k} = 0 \\
&\boxed{\phi_{a \leftarrow b,k} \propto \exp \left\{ \mu_{b,k} + \phi_{a \rightarrow b,k} (\psi(b_{k0}) - \psi(b_{k0} + b_{k1}) - \ln \epsilon) \right\}}
\end{aligned}$$

#### 4.8 Gradient with respect to $\phi_{a \rightarrow b, k}$ for nonlinks

$$\begin{aligned}
\mathcal{L}_{\phi_{a \rightarrow b, k}} &= \phi_{a \rightarrow b, k} \mu_{a, k} \\
&\quad + \phi_{a \rightarrow b, k} \phi_{a \leftarrow b, k} (\psi(b_{k1}) - \psi(b_{k0} + b_{k1}) - \ln(1 - \epsilon)) \\
&\quad - \phi_{a \rightarrow b, k} \ln \phi_{a \rightarrow b, k} \\
&= \phi_{a \rightarrow b, k} \left( \mu_{a, k} + \phi_{a \leftarrow b, k} (\psi(b_{k1}) - \psi(b_{k0} + b_{k1}) - \ln(1 - \epsilon)) - \ln \phi_{a \rightarrow b, k} \right) \\
\nabla_{\phi_{a \rightarrow b, k}} \mathcal{L}_{\phi_{a \rightarrow b, k}} &= \mu_{a, k} + \phi_{a \leftarrow b, k} (\psi(b_{k1}) - \psi(b_{k0} + b_{k1}) - \ln(1 - \epsilon)) - \ln \phi_{a \rightarrow b, k} = 0 \\
&\quad \boxed{\phi_{a \rightarrow b, k} \propto \exp \left\{ \mu_{a, k} + \phi_{a \leftarrow b, k} (\psi(b_{k1}) - \psi(b_{k0} + b_{k1}) - \ln(1 - \epsilon)) \right\}}
\end{aligned}$$

#### 4.9 Gradient with respect to $\phi_{a \leftarrow b, k}$ for nonlinks

$$\begin{aligned}
\mathcal{L}_{\phi_{a \leftarrow b, k}} &= \phi_{a \leftarrow b, k} \mu_{b, k} \\
&\quad + \phi_{a \rightarrow b, k} \phi_{a \leftarrow b, k} (\psi(b_{k1}) - \psi(b_{k0} + b_{k1}) - \ln(1 - \epsilon)) \\
&\quad - \phi_{a \leftarrow b, k} \ln \phi_{a \leftarrow b, k} \\
&= \phi_{a \leftarrow b, k} \left( \mu_{b, k} + \phi_{a \rightarrow b, k} (\psi(b_{k1}) - \psi(b_{k0} + b_{k1}) - \ln(1 - \epsilon)) - \ln \phi_{a \leftarrow b, k} \right) \\
\nabla_{\phi_{a \leftarrow b, k}} \mathcal{L}_{\phi_{a \leftarrow b, k}} &= \mu_{b, k} + \phi_{a \rightarrow b, k} (\psi(b_{k1}) - \psi(b_{k0} + b_{k1}) - \ln(1 - \epsilon)) - \ln \phi_{a \leftarrow b, k} = 0 \\
&\quad \boxed{\phi_{a \leftarrow b, k} \propto \exp \left\{ \mu_{b, k} + \phi_{a \rightarrow b, k} (\psi(b_{k1}) - \psi(b_{k0} + b_{k1}) - \ln(1 - \epsilon)) \right\}}
\end{aligned}$$

#### 4.10 Gradient with respect to $\mu_a$

$\mu_a$  and  $\Lambda_a$  are two of the scarier ones.

$$\begin{aligned}
\mathcal{L}_{\mu_{a,k}} &= -\frac{\ell}{2}[(\mu_{a,k} - m_k)^T L_{kk}(\mu_{a,k} - m_k) + \\
&\quad \sum_{b \in \text{sink}(a)} \phi_{a \rightarrow b,k} \mu_{a,k} + \\
&\quad \sum_{b \notin \text{sink}(a)} \phi_{a \rightarrow b,k} \mu_{a,k} + \\
&\quad \sum_{b \in \text{source}(a)} \phi_{b \leftarrow a,k} \mu_{a,k} + \\
&\quad \sum_{b \notin \text{source}(a)} \phi_{b \leftarrow a,k} \mu_{a,k} - \\
&\quad \sum_b \zeta_a^{-1} \sum_l \exp(\mu_{a,l} + \frac{1}{2} \text{diag}(\Lambda_a^{-1})_{ll}) \\
&=
\end{aligned}$$

$$\nabla_{\mu_{a,k}} \mathcal{L}_{\mu_{a,k}} = -\ell L_{kk}(\mu_{a,k} - m_k) + \sum_b (\phi_{a \rightarrow b,k} + \phi_{b \leftarrow a,k}) - \sum_b \zeta_a^{-1} \exp(\mu_{a,k} + \frac{1}{2} \text{diag}(\Lambda_a^{-1})_{ll})$$

Also

$$\begin{aligned}
\mathcal{L}_{\zeta_a} &= -\zeta_a^{-1} \sum_l \exp(\mu_{a,l} + \frac{1}{2} \text{diag}(\Lambda_a^{-1})_{ll}) - \ln \zeta_a \\
\nabla_{\zeta_a} \mathcal{L}_{\zeta_a} &= -\zeta_a^{-2} \sum_l \exp(\mu_{a,l} + \frac{1}{2} \text{diag}(\Lambda_a^{-1})_{ll}) - \zeta_a^{-1} = 0
\end{aligned}$$

$$= \boxed{\sum_l \exp(\mu_{a,l} + \frac{1}{2} \text{diag}(\Lambda_a^{-1})_{ll}) = \zeta_a}$$

$$\text{So } \zeta_a^{-1} \exp(\mu_{a,k} + \frac{1}{2} \text{diag}(\Lambda_a^{-1})_{ll}) = \text{softmax}(\mu_{a,k} + \frac{1}{2} \text{diag}(\Lambda_a^{-1})_{ll})$$

$$\nabla_{\mu_{a,k}} \mathcal{L}_{\mu_{a,k}} = G_{\mu_{a,k}} \propto \boxed{-\ell L_{kk}(\mu_{a,k} - m_k) + \sum_b (\phi_{a \rightarrow b,k} + \phi_{b \leftarrow a,k}) - \sum_b \text{softmax}(\cdot) (1 - \text{softmax}(\cdot))}$$

$$\nabla_{\mu_{a,k}}^2 \mathcal{L}_{\mu_{a,k}} = H_{\mu_{a,k}} \propto \boxed{-\ell L_{kk} - \sum_b \left( \text{softmax}(\cdot) - 3\text{softmax}^2(\cdot) + 2\text{softmax}^3(\cdot) \right)}$$

The newton step would look like:

$$\mu_{a,k} = \mu_{a,k} - H_{\mu_{a,k}}^{-1} G_{\mu_{a,k}}$$

#### 4.11 Gradient with respect to $\Lambda_a$

$$\begin{aligned}
\mathcal{L}_{\Lambda_{a,k}^{-1}} &= -\frac{\ell}{2} \text{tr}(L \Lambda_a^{-1}) + \frac{1}{2} \ln |\Lambda_a^{-1}| - \sum_b \zeta_a^{-1} \exp(\mu_{a,k} + \frac{1}{2} \Lambda_{a,k}^{-1}) \\
&= -\frac{\ell}{2} (L_{kk} \Lambda_{a,k}^{-1}) - \frac{1}{2} \sum_{j \neq k} \ln \Lambda_{a,j} - \frac{1}{2} \ln \Lambda_{a,k} - \sum_b \zeta_a^{-1} \exp(\mu_{a,k} + \frac{1}{2} \Lambda_{a,k}^{-1}) \\
\nabla_{\Lambda_{a,k}^{-1}} \mathcal{L}_{\Lambda_{a,k}^{-1}} &= G_{\Lambda_{a,k}^{-1}} = \boxed{-\frac{\ell}{2} L_{kk} + \frac{1}{2} \Lambda_{a,k} - \frac{1}{2} \sum_b \text{softmax}(\cdot) (1 - \text{softmax}(\cdot))}
\end{aligned}$$

□

$$\nabla_{\Lambda_{a,k}^{-1}}^2 \mathcal{L}_{\Lambda_{a,k}^{-1}} = H_{\Lambda_{a,k}^{-1}} \propto \boxed{-\frac{1}{2} \Lambda_{a,k}^{-2} - \frac{1}{4} \sum_b \left( \text{softmax}(\cdot) - 3\text{softmax}^2(\cdot) + 2\text{softmax}^3(\cdot) \right)}$$

The newton step would look like:

$$\Lambda_{a,k}^{-1} = \Lambda_{a,k}^{-1} - H_{\Lambda_{a,k}^{-1}}^{-1} G_{\Lambda_{a,k}^{-1}}$$

## 5 Other notes

### 5.1 Checknig ELBO

I need to use the training data, perhaps all the training links and same number sampled from nonlinks. Depending on the size of the network and the sparsity, this could be costly, so this only is needed every once in a while(not frequently) and only to check whether the ELBO is improving with our optimization algorithm or not.

### 5.2 some fixed ones

I should ensure that  $M$  remains  $I + N\ell L$  with  $N$  the actual size of the nodes in the training set. This also affects the update of  $m$  which is  $M^{-1}(\ell L \frac{N}{\#mbnodes} \sum_{a \in mbnodes} \mu_a)$ . Moreover  $\ell = K + N$  can be fixed in advanced, and there is no need for computation in the variational loop.

### 5.3 Newton step behavior

I should ensure if the ELBO does not improve, I should half the step in the learning rate, other wise keep it the same. In general ELBO should be in sample and not out-of-sample, so I should think whether I want this ELBO to be on a sample training or just the minibatch.

### 5.4 Initialization

I should check whether the initialization is very off or should I adopt the Gopalan's initialization algorithm.