

# Bayesian Inference for Sparse Factor Models

Yong See Foo

January 2020

## 1 Model

Suppose that we have a matrix  $\mathbf{Y} \in \mathbb{R}^{G \times N}$  of observed gene expressions, where  $G$  is the number of genes, and  $N$  is the number of individuals. We wish to model gene expression as a weighted sum of  $K$  transcription factor activities:  $y_{ij} = \sum_{k=1}^K l_{ik} f_{kj} + e_{ij}$ , where  $l_{ik}$  is the regulatory weight of factor  $k$  on gene  $i$ ,  $f_{kj}$  is the activation of factor  $k$  for individual  $j$ , and  $e_{ij}$  accounts for any corresponding residual noise. In matrix notation, the model is formulated as  $\mathbf{Y} = \mathbf{L}\mathbf{F} + \mathbf{E}$ . By assuming that the noise is independently distributed and follows a Gaussian distribution with gene-specific variance, the distribution of the gene expression  $\mathbf{Y}$  can be defined as

$$p(\mathbf{y}_i \mid \mathbf{L}, \mathbf{F}, \boldsymbol{\tau}) = \mathcal{N}(\mathbf{y}_i \mid \mathbf{F}^\top \mathbf{l}_i, \tau_i^{-1} \mathbf{I}),$$

where  $\mathbf{y}_i$  and  $\mathbf{l}_i$  are column vectors indicating the  $i$ th row of  $\mathbf{Y}$  and  $\mathbf{L}$  respectively, and  $\tau_i$  is the precision of Gaussian noise. This may also be written as  $p(\mathbf{y}_i \mid \mathbf{L}, \mathbf{F}, \boldsymbol{\tau}) = \mathcal{N}(\mathbf{y}_i \mid \mathbf{L} \mathbf{f}_i, D_{\boldsymbol{\tau}}^{-1})$ , where  $\mathbf{f}_i$  indicates the  $i$ th column of  $\mathbf{F}$ , and  $D_{\mathbf{v}} = \text{diag}(\mathbf{v})$  for any vector  $\mathbf{v}$ .

As only a small subset of genes are regulated by each transcription factor, the loading matrix  $\mathbf{L}$  is known to be sparse. This is encoded with the following prior:

$$p(l_{ik} \mid z_{ik}, \alpha_k) = \begin{cases} \delta(l_{ik}) & \text{if } z_{ik} = 0 \\ \mathcal{N}(l_{ik} \mid 0, \alpha_k^{-1}) & \text{if } z_{ik} = 1 \end{cases}$$

where  $z_{ik} = 0$  if gene  $i$  is not regulated by transcription factor  $k$ , otherwise  $l_{ik}$  follows a Gaussian distribution with factor-specific precision  $\alpha_k$ . A connectivity matrix  $\mathbf{Z}$  stores the latent binary variables  $z_{ik}$ , and we define a Bernoulli prior for each of its elements:

$$p(z_{ik}) = \text{Bernoulli}(z_{ik} \mid \pi_k),$$

where  $\pi_k$  are hyperparameters which control the sparsity of each factor.

To avoid identifiability issues caused by scaling, we define a unit Gaussian prior distribution for the factor matrix  $\mathbf{F}$ :

$$p(\mathbf{f}_j) = \mathcal{N}(\mathbf{f}_j \mid \mathbf{0}, \mathbf{I}).$$

Lastly, a gamma prior is defined for each of the precision parameters:

$$\begin{aligned} p(\tau_i) &= \Gamma(\tau_i \mid a_\tau, b_\tau) \\ p(\alpha_k) &= \Gamma(\alpha_k \mid a_\alpha, b_\alpha), \end{aligned}$$

where  $a_\tau, b_\tau, a_\alpha, b_\alpha$  are hyperparameters to be specified.

## 2 MCMC

### 2.1 Full conditionals for Gibbs sampling

We need to ensure that  $l_{ik} = 0$  whenever  $z_{ik} = 0$ . This can be achieved by introducing modifications to the conditional distribution of  $\mathbf{Y}$ . We have

$$\begin{aligned}
p(\mathbf{l}_i, \mathbf{z}_i \mid \mathbf{Y}, \mathbf{F}, \boldsymbol{\tau}, \boldsymbol{\alpha}) &\propto \prod_{k: z_{ik}=1} \pi_k \sqrt{\frac{\alpha_k}{2\pi}} \times \prod_{k: z_{ik}=0} (1 - \pi_k) \delta(l_{ik}) \\
&\quad \times \exp \left\{ -\frac{\tau_i}{2} \left( \mathbf{y}_i - [\mathbf{F}]_{\mathbf{z}_i}^\top [\mathbf{l}_i]_{\mathbf{z}_i} \right)^\top \left( \mathbf{y}_i - [\mathbf{F}]_{\mathbf{z}_i}^\top [\mathbf{l}_i]_{\mathbf{z}_i} \right) - \frac{1}{2} [\mathbf{l}_i]_{\mathbf{z}_i}^\top [D\boldsymbol{\alpha}]_{\mathbf{z}_i} [\mathbf{l}_i]_{\mathbf{z}_i} \right\} \\
&\propto \prod_{k: z_{ik}=1} \pi_k \sqrt{\frac{\alpha_k}{2\pi}} \times \prod_{k: z_{ik}=0} (1 - \pi_k) \delta(l_{ik}) \\
&\quad \times \exp \left\{ -\frac{1}{2} ([\mathbf{l}]_{\mathbf{z}_i} - \boldsymbol{\mu}_{\mathbf{l}_i})^\top \Sigma_{\mathbf{l}_i}^{-1} ([\mathbf{l}]_{\mathbf{z}_i} - \boldsymbol{\mu}_{\mathbf{l}_i}) + \frac{1}{2} \boldsymbol{\mu}_{\mathbf{l}_i}^\top \Sigma_{\mathbf{l}_i}^{-1} \boldsymbol{\mu}_{\mathbf{l}_i} \right\}
\end{aligned} \tag{1}$$

where

$$\begin{aligned}
[\mathbf{F}]_{\mathbf{z}_i} &= \text{matrix consisting of rows of } \mathbf{F} \text{ whose corresponding entries of } \mathbf{z}_i \text{ are equal to 1} \\
[\mathbf{l}_i]_{\mathbf{z}_i} &= \text{vector consisting of entries of } \mathbf{l}_i \text{ whose corresponding entries of } \mathbf{z}_i \text{ are equal to 1} \\
[D\boldsymbol{\alpha}]_{\mathbf{z}_i} &= \text{matrix consisting of rows of } D\boldsymbol{\alpha} \text{ whose corresponding entries of } \mathbf{z}_i \text{ are equal to 1} \\
\Sigma_{\mathbf{l}_i} &= \left( \tau_i [\mathbf{F}]_{\mathbf{z}_i} [\mathbf{F}]_{\mathbf{z}_i}^\top + [D\boldsymbol{\alpha}]_{\mathbf{z}_i} \right)^{-1} \\
\boldsymbol{\mu}_{\mathbf{l}_i} &= \tau_i \Sigma_{\mathbf{l}_i} [\mathbf{F}]_{\mathbf{z}_i}^\top \mathbf{y}_i,
\end{aligned}$$

and hence obtain the full conditional distribution of  $\mathbf{l}_i$ :

$$p([\mathbf{l}_i]_{\mathbf{z}_i} \mid \mathbf{Y}, \mathbf{F}, \mathbf{Z}, \boldsymbol{\tau}, \boldsymbol{\alpha}) = \mathcal{N}([\mathbf{l}_i]_{\mathbf{z}_i} \mid \boldsymbol{\mu}_{\mathbf{l}_i}, \Sigma_{\mathbf{l}_i}) \times \prod_{k: z_{ik}=0} \delta(l_{ik}). \tag{2}$$

Marginalising out  $\mathbf{l}_i$  from Equation 1 gives a conditional distribution of  $z_{ik}$ :

$$p(z_{ik} \mid \mathbf{Y}, \mathbf{F}, \mathbf{Z}_{-ik}, \boldsymbol{\tau}, \boldsymbol{\alpha}) \propto \left( \frac{\alpha_k}{2\pi} \right)^{\frac{z_{ik}}{2}} \det |\Sigma_{\mathbf{l}_i}|^{\frac{1}{2}} \exp \left\{ \frac{1}{2} \boldsymbol{\mu}_{\mathbf{l}_i}^\top \Sigma_{\mathbf{l}_i}^{-1} \boldsymbol{\mu}_{\mathbf{l}_i} \right\} \pi_k^{z_{ik}} (1 - \pi_k)^{1-z_{ik}}. \tag{3}$$

We also have

$$\begin{aligned}
p(\mathbf{f}_j \mid \mathbf{Y}, \mathbf{L}, \mathbf{Z}, \boldsymbol{\tau}, \boldsymbol{\alpha}) &\propto \exp \left\{ -\frac{1}{2} (\mathbf{y}_j - \mathbf{L} \mathbf{f}_j)^\top D\boldsymbol{\tau} (\mathbf{y}_j - \mathbf{L} \mathbf{f}_j) - \frac{1}{2} \mathbf{f}_j^\top \mathbf{f}_j \right\} \\
&\propto \exp \left\{ -\frac{1}{2} (\mathbf{f}_j - \boldsymbol{\mu}_{\mathbf{f}_j})^\top \Sigma_{\mathbf{f}_j}^{-1} (\mathbf{f}_j - \boldsymbol{\mu}_{\mathbf{f}_j}) \right\}
\end{aligned}$$

where

$$\begin{aligned}
\Sigma_{\mathbf{f}_j} &= \left( \mathbf{L}^\top D\boldsymbol{\tau} \mathbf{L} + \mathbf{I} \right)^{-1} \\
\boldsymbol{\mu}_{\mathbf{f}_j} &= \Sigma_{\mathbf{f}_j} \mathbf{L}^\top D\boldsymbol{\tau} \mathbf{y}_j,
\end{aligned}$$

thus arriving at the full conditional distribution of  $\mathbf{f}_j$ :

$$p(\mathbf{f}_j \mid \mathbf{Y}, \mathbf{L}, \mathbf{Z}, \boldsymbol{\tau}, \boldsymbol{\alpha}) = \mathcal{N}(\mathbf{f}_j \mid \boldsymbol{\mu}_{\mathbf{f}_j}, \Sigma_{\mathbf{f}_j}). \tag{4}$$

Lastly, we have the full conditional distribution of  $\tau_i$ :

$$p(\tau_i \mid \mathbf{Y}, \mathbf{L}, \mathbf{F}, \mathbf{Z}, \boldsymbol{\alpha}) = \Gamma \left( \tau_i \mid a_\tau + \frac{N}{2}, b_\tau + \frac{1}{2} \left( \mathbf{y}_i - [\mathbf{F}]_{\mathbf{z}_i}^\top [\mathbf{l}_i]_{\mathbf{z}_i} \right)^\top \left( \mathbf{y}_i - [\mathbf{F}]_{\mathbf{z}_i}^\top [\mathbf{l}_i]_{\mathbf{z}_i} \right) \right), \tag{5}$$

and the full conditional distribution of  $\alpha_k$ :

$$p(\alpha_k \mid \mathbf{Y}, \mathbf{L}, \mathbf{F}, \mathbf{Z}, \boldsymbol{\tau}) = \Gamma\left(\alpha_k \mid a_\alpha + \frac{1}{2} \sum_{i=1}^G z_{ik}, b_\alpha + \frac{1}{2} \sum_{i: z_{ik}=1} l_{ik}^2\right). \quad (6)$$

## 2.2 Relabelling

A relabelling algorithm is used to deal with signflip symmetries and label switching.

## 3 Variational Inference

### 3.1 Mean-field approximation

Use the variational factorisation

$$q(\mathbf{L}, \mathbf{F}, \mathbf{Z}, \boldsymbol{\tau}, \boldsymbol{\alpha}) = \prod_{i=1}^G \left[ \prod_{k=1}^K q(l_{ik} \mid z_{ik}) q(z_{ik}) \right] q(\boldsymbol{\tau}_i) \times \prod_{j=1}^N q(\mathbf{f}_j) \times \prod_{k=1}^K q(\alpha_k) \quad (7)$$

as an approximation to the posterior distribution, where

$$\begin{aligned} q(l_{ik} \mid z_{ik}) &= \mathcal{N}(l_{ik} \mid \mu_{l_{ik}}, \sigma_{l_{ik}}^2)^{z_{ik}} \times \delta(l_{ik})^{1-z_{ik}} \\ q(z_{ik}) &= \text{Bernoulli}(z_{ik} \mid \eta_{ik}) \\ q(\mathbf{f}_j) &= \mathcal{N}(\mathbf{f}_j \mid \boldsymbol{\mu}_{\mathbf{f}_j}, \Sigma_{\mathbf{f}_j}) \\ q(\boldsymbol{\tau}_i) &= \Gamma(\boldsymbol{\tau}_i \mid \hat{a}_{\boldsymbol{\tau}_i}, \hat{b}_{\boldsymbol{\tau}_i}) \\ q(\alpha_k) &= \Gamma(\alpha_k \mid \hat{a}_{\alpha_k}, \hat{b}_{\alpha_k}). \end{aligned}$$

Coordinate ascent for  $\mathbf{l}_i$  and  $\mathbf{z}_i$  gives

$$\begin{aligned} q^*(l_{ik}, z_{ik}) &\propto \exp \left\{ \mathbb{E}_{\mathbf{l}_{-ik}, \mathbf{F}, \mathbf{z}_{-ik}, \boldsymbol{\tau}_i, \boldsymbol{\alpha}} [\log p(l_{ik}, z_{ik} \mid \mathbf{Y}, \mathbf{l}_{-ik}, \mathbf{F}, \mathbf{z}_{-ik}, \boldsymbol{\tau}, \boldsymbol{\alpha})] \right\} \\ &\propto \exp \left\{ \frac{z_{ik}}{2} \mathbb{E}_{\alpha_k} \left[ \log \frac{\alpha_k}{2\pi} \right] - \mathbb{E}_{\mathbf{l}_{-ik}, \mathbf{F}, \mathbf{z}_{-ik}, \boldsymbol{\tau}_i} \left[ \frac{\tau_i}{2} \left( \mathbf{y}_i - [\mathbf{F}]_{\mathbf{z}_i}^\top [\mathbf{l}]_{\mathbf{z}_i} \right)^\top \left( \mathbf{y}_i - [\mathbf{F}]_{\mathbf{z}_i}^\top [\mathbf{l}]_{\mathbf{z}_i} \right) \right] \right. \\ &\quad \left. - \frac{z_{ik}}{2} \mathbb{E}_{\boldsymbol{\alpha}} [\alpha_k l_{ik}^2] \right\} \times \prod_{k=1}^K \pi_k^{z_{ik}} (1 - \pi_k)^{1-z_{ik}}, \\ &\propto \exp \left\{ -\frac{\hat{a}_{\tau_i}}{2\hat{b}_{\tau_i}} \mathbb{E}_{\mathbf{l}_{-ik}, \mathbf{F}, \mathbf{z}_{-ik}} \left[ -2z_{ik} \mathbf{y}_i^\top \mathbf{f}_k l_{ik} + 2z_{ik} \mathbf{f}_k^\top \sum_{k' \neq k} z_{ik'} \mathbf{f}_{k'} l_{ik'} l_{ik} + z_{ik} \mathbf{f}_k^\top \mathbf{f}_k l_{ik}^2 \right] \right. \\ &\quad \left. + \frac{z_{ik}}{2} \left( \psi(\hat{a}_{\alpha_k}) - \log 2\pi \hat{b}_{\alpha_k} - \frac{\hat{a}_{\alpha_k}}{\hat{b}_{\alpha_k}} l_{ik}^2 \right) \right\} \times \prod_{k=1}^K \pi_k^{z_{ik}} (1 - \pi_k)^{1-z_{ik}} \\ &\propto \exp \left\{ -\frac{\hat{a}_{\tau_i}}{2\hat{b}_{\tau_i}} \left[ -2z_{ik} \mathbf{y}_i^\top \boldsymbol{\mu}_{\mathbf{f}_k} l_{ik} + 2z_{ik} \boldsymbol{\mu}_{\mathbf{f}_k}^\top \sum_{k' \neq k} \eta_{ik'} \boldsymbol{\mu}_{\mathbf{f}_{k'}} \boldsymbol{\mu}_{\mathbf{l}_{ik'}} l_{ik} \right. \right. \\ &\quad \left. \left. + z_{ik} \sum_{j=1}^N \left( [\Sigma_{\mathbf{f}_j}]_{kk} + [\boldsymbol{\mu}_{\mathbf{f}_j}]_k^2 \right) l_{ik}^2 \right] + \frac{z_{ik}}{2} \left( \psi(\hat{a}_{\alpha_k}) - \log 2\pi \hat{b}_{\alpha_k} - \frac{\hat{a}_{\alpha_k}}{\hat{b}_{\alpha_k}} l_{ik}^2 \right) \right\} \\ &\quad \times \pi_k^{z_{ik}} (1 - \pi_k)^{1-z_{ik}}, \end{aligned}$$

which corresponds to the updates

$$\sigma_{l_{ik}}^{2*} = \left( \frac{\hat{a}_{\tau_i}}{\hat{b}_{\tau_i}} \sum_{j=1}^N \left( [\Sigma \mathbf{f}_j]_{kk} + [\boldsymbol{\mu} \mathbf{f}_j]_k^2 \right) + \frac{\hat{a}_{\alpha_k}}{\hat{b}_{\alpha_k}} \right)^{-1} \quad (8)$$

$$\mu_{l_{ik}}^* = \frac{\hat{a}_{\tau_i}}{\hat{b}_{\tau_i}} \sigma_{l_{ik}}^{2*} \boldsymbol{\mu}_{\mathbf{f}_k}^\top \left( \mathbf{y}_i - \sum_{k' \neq k} \eta_{ik'} \boldsymbol{\mu}_{\mathbf{f}_{k'}} \mu_{l_{ik'}} \right)$$

$$q(z_{ik}) \propto \exp \left\{ \frac{z_{ik}}{2} \left( \psi(\hat{a}_{\alpha_k}) - \log 2\pi \hat{b}_{\alpha_k} + \frac{\mu_{l_{ik}}^{2*}}{\sigma_{l_{ik}}^{2*}} \right) \right\} \left( \sqrt{\sigma_{l_{ik}}^{2*}} \pi_k \right)^{z_{ik}} (1 - \pi_k)^{1-z_{ik}}. \quad (9)$$

Coordinate ascent for  $\mathbf{f}_j$  gives

$$\begin{aligned} q^*(\mathbf{f}_j) &\propto \exp \{ \mathbb{E}_{\mathbf{L}, \mathbf{Z}, \boldsymbol{\tau}} [\log p(\mathbf{f}_j \mid \mathbf{Y}, \mathbf{L}, \mathbf{Z}, \boldsymbol{\tau}, \boldsymbol{\alpha})] \} \\ &\propto \exp \left\{ \mathbb{E}_{\mathbf{L}, \mathbf{L}, \mathbf{Z}, \boldsymbol{\tau}} \left[ -\frac{1}{2} (\mathbf{y}_j - \mathbf{L} \mathbf{f}_j)^\top D_{\boldsymbol{\tau}} (\mathbf{y}_j - \mathbf{L} \mathbf{f}_j) \right] - \frac{1}{2} \mathbf{f}_j^\top \mathbf{f}_j \right\} \\ &\propto \exp \left\{ \mathbf{y}_j^\top D_{\bar{\boldsymbol{\tau}}} \bar{\mathbf{L}} \mathbf{f}_j - \frac{1}{2} \mathbf{f}_j^\top \bar{\mathbf{L}}^\top D_{\boldsymbol{\tau}} \bar{\mathbf{L}} \mathbf{f}_j - \frac{1}{2} \mathbf{f}_j^\top \mathbf{f}_j \right\} \end{aligned}$$

where

$$\begin{aligned} D_{\bar{\boldsymbol{\tau}}} &= \text{diag} \left( \left\{ \frac{\hat{a}_{\tau_i}}{\hat{b}_{\tau_i}} \right\}_{i=1}^G \right) \\ [\bar{\mathbf{L}}]_{ik} &= \eta_{ik} \mu_{l_{ik}} \\ [\bar{\mathbf{L}}^\top D_{\boldsymbol{\tau}} \bar{\mathbf{L}}]_{kk'} &= \sum_{i=1}^G \frac{\hat{a}_{\tau_i}}{\hat{b}_{\tau_i}} \eta_{ik} \eta_{ik'}^{1-\delta_{kk'}} (\delta_{kk'} \sigma_{l_{ik}}^2 + \mu_{l_{ik}} \mu_{l_{ik'}}), \end{aligned}$$

which corresponds to the updates

$$\begin{aligned} \Sigma_{\mathbf{f}_j}^* &= \left( \bar{\mathbf{L}}^\top D_{\boldsymbol{\tau}} \bar{\mathbf{L}} + \mathbf{I} \right)^{-1} \\ \boldsymbol{\mu}_{\mathbf{f}_j}^* &= \Sigma_{\mathbf{f}_j}^* \bar{\mathbf{L}}^\top D_{\bar{\boldsymbol{\tau}}} \mathbf{y}_j. \end{aligned} \quad (10)$$

Coordinate ascent for  $\tau_i$  gives

$$\begin{aligned} q^*(\tau_i) &\propto \exp \{ \mathbb{E}_{\mathbf{L}, \mathbf{F}, \mathbf{Z}} [\log p(\tau_i \mid \mathbf{Y}, \mathbf{L}, \mathbf{F}, \mathbf{Z}, \boldsymbol{\alpha})] \} \\ &\propto \exp \left\{ \left( a_\tau + \frac{N}{2} \right) \log \tau_i - b_\tau \tau_i - \frac{\tau_i}{2} \mathbb{E}_{\mathbf{L}, \mathbf{F}, \mathbf{Z}} \left[ \left( \mathbf{y}_i - \mathbf{F}^\top \mathbf{l}_i \right)^\top \left( \mathbf{y}_i - \mathbf{F}^\top \mathbf{l}_i \right) \right] \right\} \\ &\propto \exp \left\{ \left( a_\tau + \frac{N}{2} \right) \log \tau_i - \left( b_\tau + \frac{1}{2} \left( \mathbf{y}_i^\top \mathbf{y}_i - 2 \bar{\mathbf{l}}_i^\top \bar{\mathbf{F}} \mathbf{y}_i + \bar{\mathbf{l}}_i^\top \bar{\mathbf{F}} \bar{\mathbf{F}}^\top \bar{\mathbf{l}}_i \right) \right) \tau_i \right\} \end{aligned}$$

where

$$\begin{aligned} \bar{\mathbf{l}}_i &= \{ \eta_{ik} \mu_{l_{ik}} \}_{k=1}^K \\ [\bar{\mathbf{F}}]_{kj} &= \mu_{f_{kj}} \\ \bar{\mathbf{l}}_i^\top \bar{\mathbf{F}} \bar{\mathbf{F}}^\top \bar{\mathbf{l}}_i &= \sum_{k=1}^K \sum_{k'=1}^K \left( \eta_{ik} \eta_{ik'}^{1-\delta_{kk'}} (\delta_{kk'} \sigma_{l_{ik}}^2 + \mu_{l_{ik}} \mu_{l_{ik'}}) \sum_{j=1}^N \left( [\Sigma \mathbf{f}_j]_{kk'} + [\boldsymbol{\mu} \mathbf{f}_j]_k [\boldsymbol{\mu} \mathbf{f}_j]_{k'} \right) \right), \end{aligned}$$

which corresponds to the updates

$$\begin{aligned}\hat{a}_{\tau_i}^* &= a_\tau + \frac{N}{2} \\ \hat{b}_{\tau_i}^* &= b_\tau + \frac{1}{2} \left( \mathbf{y}_i^\top \mathbf{y}_i - 2\bar{\mathbf{l}}_i^\top \bar{\mathbf{F}} \mathbf{y}_i + \overline{\mathbf{l}_i^\top \mathbf{F} \mathbf{F}^\top \mathbf{l}_i} \right).\end{aligned}\tag{11}$$

Coordinate ascent for  $\alpha_k$  gives

$$\begin{aligned}q^*(\alpha_k) &\propto \exp \{ \mathbb{E}_{\mathbf{L}, \mathbf{Z}} [\log p(\alpha_k \mid \mathbf{Y}, \mathbf{L}, \mathbf{F}, \mathbf{Z}, \boldsymbol{\tau})] \} \\ &\propto \exp \left\{ \left( a_\alpha + \frac{1}{2} \mathbb{E}_{\mathbf{Z}} \left[ \sum_{i=1}^G z_{ik} \right] \right) \log \alpha_k - b_\alpha \alpha_k - \frac{\alpha_k}{2} \mathbb{E}_{\mathbf{L}, \mathbf{Z}} \left[ \sum_{i: z_{ik}=1} l_{ik}^2 \right] \right\}\end{aligned}$$

which corresponds to the updates

$$\begin{aligned}\hat{a}_{\alpha_k}^* &= a_\alpha + \frac{1}{2} \sum_{i=1}^G \eta_{ik} \\ \hat{b}_{\alpha_k}^* &= b_\alpha + \frac{1}{2} \sum_{i=1}^G \eta_{ik} (\sigma_{l_{ik}}^2 + \mu_{l_{ik}}^2).\end{aligned}\tag{12}$$

## 4 Unused attempts

### 4.1 Capturing dependency within $\mathbf{l}_i$

To capture dependency within  $\mathbf{l}_i$ , the variational factorisation is modified:

$$q(\mathbf{L}, \mathbf{F}, \mathbf{Z}, \boldsymbol{\tau}, \boldsymbol{\alpha}) = \prod_{i=1}^G q(\mathbf{l}_i, \mathbf{z}_i) q(\tau_i) \times \prod_{j=1}^N q(\mathbf{f}_j) \times \prod_{k=1}^K q(\alpha_k)\tag{13}$$

where

$$\begin{aligned}q(\mathbf{l}_i, \mathbf{z}_i) &= \mathcal{N}([\mathbf{l}_i]_{\mathbf{z}_i} \mid \boldsymbol{\mu}_{\mathbf{l}_i}, \Sigma_{\mathbf{l}_i}) \times \prod_{k: z_{ik}=0} \delta(l_{ik}) \times q(\mathbf{z}_i) \\ q(\mathbf{f}_j) &= \mathcal{N}(\mathbf{f}_j \mid \boldsymbol{\mu}_{\mathbf{f}_j}, \Sigma_{\mathbf{f}_j}) \\ q(\tau_i) &= \Gamma(\tau_i \mid \hat{a}_{\tau_i}, \hat{b}_{\tau_i}) \\ q(\alpha_k) &= \Gamma(\alpha_k \mid \hat{a}_{\alpha_k}, \hat{b}_{\alpha_k}).\end{aligned}$$

Note that the dimensions of  $\boldsymbol{\mu}_{\mathbf{l}_i}$  and  $\Sigma_{\mathbf{l}_i}$  depend on  $\mathbf{z}_i$ , so there is no sensible way to numerically record its value. This will prove to be problematic.

Coordinate ascent for  $\mathbf{l}_i$  and  $\mathbf{z}_i$  gives

$$\begin{aligned}
q^*(\mathbf{l}_i, \mathbf{z}_i) &\propto \exp \{ \mathbb{E}_{\mathbf{F}, \tau_i, \alpha} [\log p(\mathbf{l}_i, \mathbf{z}_i \mid \mathbf{Y}, \mathbf{F}, \tau, \alpha)] \} \\
&\propto \exp \left\{ \frac{1}{2} \sum_{k: z_{ik}=1} \mathbb{E}_{\alpha_k} \left[ \log \frac{\alpha_k}{2\pi} \right] - \mathbb{E}_{\mathbf{F}, \tau_i} \left[ \frac{\tau_i}{2} \left( \mathbf{y}_i - [\mathbf{F}]_{\mathbf{z}_i}^\top [\mathbf{l}_i]_{\mathbf{z}_i} \right)^\top \left( \mathbf{y}_i - [\mathbf{F}]_{\mathbf{z}_i}^\top [\mathbf{l}_i]_{\mathbf{z}_i} \right) \right] \right. \\
&\quad \left. - \frac{1}{2} \mathbb{E}_{\alpha} \left[ [\mathbf{l}_i]_{\mathbf{z}_i}^\top [D\alpha]_{\mathbf{z}_i} [\mathbf{l}_i]_{\mathbf{z}_i} \right] \right\} \times \prod_{k=1}^K \pi_k^{z_{ik}} (1 - \pi_k)^{1-z_{ik}} \\
&\propto \exp \left\{ \frac{1}{2} \sum_{k: z_{ik}=1} (\overline{\log \alpha_k} - \log 2\pi) - \frac{\hat{a}_{\tau_i}}{2\hat{b}_{\tau_i}} \mathbb{E}_{\mathbf{F}} \left[ -2\mathbf{y}_i^\top [\mathbf{F}]_{\mathbf{z}_i}^\top [\mathbf{l}_i]_{\mathbf{z}_i} + [\mathbf{l}_i]_{\mathbf{z}_i}^\top [\mathbf{F}]_{\mathbf{z}_i} [\mathbf{F}]_{\mathbf{z}_i}^\top [\mathbf{l}_i]_{\mathbf{z}_i} \right] \right. \\
&\quad \left. - \frac{1}{2} [\mathbf{l}_i]_{\mathbf{z}_i}^\top [D\bar{\alpha}]_{\mathbf{z}_i} [\mathbf{l}_i]_{\mathbf{z}_i} \right\} \times \prod_{k=1}^K \pi_k^{z_{ik}} (1 - \pi_k)^{1-z_{ik}} \\
&\propto \exp \left\{ \frac{1}{2} \sum_{k: z_{ik}=1} (\overline{\log \alpha_k} - \log 2\pi) + \frac{\hat{a}_{\tau_i}}{\hat{b}_{\tau_i}} \sum_{k: z_{ik}=1} l_{ik} \sum_{j=1}^N y_{ij} \overline{f_{kj}} \right. \\
&\quad \left. - \frac{\hat{a}_{\tau_i}}{2\hat{b}_{\tau_i}} \sum_{k: z_{ik}=1} \sum_{k': z_{ik'}=1} l_{ik} l_{ik'} \sum_{j=1}^N \overline{f_{kj} f_{k'j}} - \frac{1}{2} [\mathbf{l}_i]_{\mathbf{z}_i}^\top [D\bar{\alpha}]_{\mathbf{z}_i} [\mathbf{l}_i]_{\mathbf{z}_i} \right\} \\
&\quad \times \prod_{k=1}^K \pi_k^{z_{ik}} (1 - \pi_k)^{1-z_{ik}},
\end{aligned}$$

where

$$\begin{aligned}
\overline{\log \alpha_k} &= \psi(\hat{a}_{\alpha_k}) - \log \hat{b}_{\alpha_k} \\
D\bar{\alpha} &= \text{diag} \left( \left\{ \frac{\hat{a}_{\alpha_k}}{\hat{b}_{\alpha_k}} \right\}_{k=1}^K \right) \\
\overline{f_{kj}} &= [\boldsymbol{\mu}_{f_j}]_k \\
\overline{f_{kj} f_{k'j}} &= [\Sigma_{f_j}]_{kk'} + [\boldsymbol{\mu}_{f_j}]_k [\boldsymbol{\mu}_{f_j}]_{k'}.
\end{aligned}$$

The variational parameters corresponding to  $[\mathbf{l}_i]_{\mathbf{z}_i}$  are thus updated to

$$\begin{aligned}
\Sigma_{\mathbf{l}_i}^* &= \left( \frac{\hat{a}_{\tau_i}}{\hat{b}_{\tau_i}} \sum_{j=1}^N \left( [\Sigma_{f_j}]_{\mathbf{z}_i, \mathbf{z}_i} + [\boldsymbol{\mu}_{f_j}]_{\mathbf{z}_i} [\boldsymbol{\mu}_{f_j}]_{\mathbf{z}_i}^\top \right) + [D\bar{\alpha}]_{\mathbf{z}_i} \right)^{-1} \\
\boldsymbol{\mu}_{\mathbf{l}_i}^* &= \frac{\hat{a}_{\tau_i}}{\hat{b}_{\tau_i}} \Sigma_{\mathbf{l}_i}^* \sum_{j=1}^N y_{ij} [\boldsymbol{\mu}_{f_j}]_{\mathbf{z}_i}
\end{aligned} \tag{14}$$

where  $[\Sigma_{f_j}]_{\mathbf{z}_i, \mathbf{z}_i}$  is a principal minor of  $\Sigma_{f_j}$  whose rows and columns indices correspond to the entries of  $\mathbf{l}_i$  in  $\mathbf{z}_i$ , and  $[\boldsymbol{\mu}_{f_j}]_{\mathbf{z}_i}$  is a vector consisting of entries of  $\boldsymbol{\mu}_{f_j}$  whose corresponding entries of  $\mathbf{z}_i$  are equal to 1.

Marginalising out  $\mathbf{l}_i$  then gives

$$q^*(\mathbf{z}_i) \propto \exp \left\{ \frac{1}{2} \sum_{k: z_{ik}=1} (\overline{\log \alpha_k} - \log 2\pi) + \frac{1}{2} \boldsymbol{\mu}_{\mathbf{l}_i}^\top \Sigma_{\mathbf{l}_i}^{-1} \boldsymbol{\mu}_{\mathbf{l}_i} \right\} \det |\Sigma_{\mathbf{l}_i}|^{\frac{1}{2}} \prod_{k=1}^K \pi_k^{z_{ik}} (1 - \pi_k)^{1-z_{ik}}. \tag{15}$$

Unfortunately, this distribution is intractable to obtain. We instead approximate  $q^*(\mathbf{z}_i)$  with

$$\hat{q}(\mathbf{z}_i) = \prod_{k=1}^K \hat{q}(z_{ik}) = \prod_{k=1}^K \text{Bernoulli}(z_{ik} \mid \gamma_{ik}). \quad (16)$$

Methods for estimating  $\gamma_{ik}$  will be addressed in Section 4.2.

For the remaining variational parameters, computations similar to the previous section give the following updates:

$$\Sigma_{\mathbf{f}_j}^* = \left( \overline{\mathbf{L}^\top D_\tau \mathbf{L}} + \mathbf{I} \right)^{-1} \quad (17)$$

$$\boldsymbol{\mu}_{\mathbf{f}_j}^* = \Sigma_{\mathbf{f}_j}^* \overline{\mathbf{L}}^\top D_\tau \mathbf{y}_j$$

$$\hat{a}_{\tau_i}^* = a_\tau + \frac{N}{2} \quad (18)$$

$$\hat{b}_{\tau_i}^* = b_\tau + \frac{1}{2} \left( \mathbf{y}_i^\top \mathbf{y}_i - 2 \overline{\mathbf{l}_i}^\top \overline{\mathbf{F}} \mathbf{y}_i + \overline{\mathbf{l}_i}^\top \mathbf{F} \mathbf{F}^\top \mathbf{l}_i \right)$$

$$\hat{a}_{\alpha_k}^* = a_\alpha + \frac{1}{2} \sum_{i=1}^G \eta_{ik} \quad (19)$$

$$\hat{b}_{\alpha_k}^* = b_\alpha + \frac{1}{2} \sum_{i=1}^G \eta_{ik} (\sigma_{l_{ik}}^2 + \mu_{l_{ik}}^2).$$

However, these expectations depend on  $\mu_{l_i}$  and  $\Sigma_{l_i}$ , which we do not have values of due to their varying dimensions.

## 4.2 Estimating $\gamma_{ik}$

Let  $B = \{0, 1\}^K$ , the set of binary vectors of size  $K$ . Through this section,  $q^*$  refers to the unnormalised density stated in Equation (15) for some fixed  $i$ . The direct approach of estimating  $\gamma_{ik}$  is to compute

$$\gamma_{ik}^* = \sum_{\substack{\boldsymbol{\zeta} \in B \\ \zeta_k=1}} q^*(\boldsymbol{\zeta}) \bigg/ \sum_{\boldsymbol{\zeta} \in B} q^*(\boldsymbol{\zeta}). \quad (20)$$

However, this would take  $2^K$  matrix inversions, which is infeasible for large  $K$ . Instead, we seek to estimate  $\gamma_{ik}$  (for all  $k$  with some fixed  $i$ ) using values of  $q^*(\boldsymbol{\zeta})$  for  $\boldsymbol{\zeta} \in \mathbf{B}_T$ , where  $\mathbf{B}_T$  is a random subset of  $B$  of size  $T$ .

Define

$$g_{ik}(z) = \sum_{\substack{\boldsymbol{\zeta} \in \mathbf{B}_T \\ \zeta_k=z}} q^*(\boldsymbol{\zeta}) \bigg/ |\{\boldsymbol{\zeta} \in \mathbf{B}_T : \zeta_k = z\}| \quad \text{for } z = 0, 1.$$

We then have an aggregation-based method for estimating  $\gamma_{ik}$ :

$$\gamma_{ik}^* = \frac{g_{ik}(1)}{g_{ik}(0) + g_{ik}(1)}. \quad (21)$$

Another *ad hoc* method is to use the independence assumption in Equation (16). Let

$$\gamma_{ik} = \frac{1}{1 + \exp(-u_{ik})},$$

the approximation

$$q^*(\zeta) \approx \prod_{k=1}^K \gamma_{ik}^{\zeta_k} (1 - \gamma_{ik})^{1-\zeta_k}$$

is then equivalent to

$$q^*(\zeta) \approx \prod_{k=1}^K \frac{\exp(u_{ik}\zeta_k)}{1 + \exp(u_{ik})}.$$

Taking logs of both sides then gives a regression problem:

$$\log q^*(\zeta) \approx u_0 + \sum_{k=1}^K u_{ik}\zeta_k \tag{22}$$

where  $u_0$  is some constant. Each element  $\zeta$  of  $\mathbf{B}_T$  and its corresponding value of  $q(\zeta)$  serves as a data point to be used for regression.

One last approach is to apply coordinate ascent to only a random subset of  $\mathbf{l}_i$  and  $\mathbf{z}_i$  during each iteration. If the subset is small enough, the direct approach found in Equation 20 can then be feasible. This idea is motivated by stochastic variational inference, but does not share the same theoretical guarantees.