

# 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

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

where  $\mathbf{y}_{i\cdot}$  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  $p(\mathbf{y}_{i\cdot} \mid \mathbf{L}, \mathbf{F}, \boldsymbol{\tau}) = \mathcal{N}(\mathbf{y}_{i\cdot} \mid \mathbf{L} \mathbf{f}_{\cdot j}, D_{\boldsymbol{\tau}}^{-1})$ , where  $\mathbf{y}_{\cdot j}$  indicates the  $j$ th column of  $\mathbf{Y}$ , 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:

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

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  $\mathbf{F}$ :

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

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 Collapsed Gibbs sampling

We use collapsed Gibbs sampling to simulate the posterior, where the regulatory weights  $\mathbf{L}$  are marginalised out when computing the conditional distribution of the connectivity matrix  $\mathbf{Z}$ . This results in a sampler that is more efficient than a vanilla Gibbs sampler, as the autocorrelation between samples of  $\mathbf{Z}$  is reduced. 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

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

where

$[\mathbf{F}]_{\mathbf{z}_{i\cdot}}$  = matrix consisting of rows of  $\mathbf{F}$  whose corresponding entries of  $\mathbf{z}_{i\cdot}$  are equal to 1

$[\mathbf{l}_{i\cdot}]_{\mathbf{z}_{i\cdot}}$  = vector consisting of entries of  $\mathbf{l}_{i\cdot}$  whose corresponding entries of  $\mathbf{z}_{i\cdot}$  are equal to 1

$[D\boldsymbol{\alpha}]_{\mathbf{z}_{i\cdot}}$  = matrix consisting of rows of  $D\boldsymbol{\alpha}$  whose corresponding entries of  $\mathbf{z}_{i\cdot}$  are equal to 1

$$\Sigma_{\mathbf{l}_{i\cdot}} = \left( \tau_i [\mathbf{F}]_{\mathbf{z}_{i\cdot}} [\mathbf{F}]_{\mathbf{z}_{i\cdot}}^\top + [D\boldsymbol{\alpha}]_{\mathbf{z}_{i\cdot}} \right)^{-1}$$

$$\boldsymbol{\mu}_{\mathbf{l}_{i\cdot}} = \tau_i \Sigma_{\mathbf{l}_{i\cdot}} [\mathbf{F}]_{\mathbf{z}_{i\cdot}}^\top \mathbf{y}_{i\cdot},$$

and hence obtain the full conditional distribution of  $\mathbf{l}_{i\cdot}$ :

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

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

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

We also have

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

where

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

thus arriving at the full conditional distribution of  $\mathbf{f}_{\cdot j}$ :

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

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

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

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 Identifiability issues

As the priors are exchangeable, this results the model being non-identifiable. Given a mode of the posterior distribution, if the latent factors are permuted, or if the sign of the entries corresponding to some factor are all switched, one will obtain another equivalent mode. These symmetries result in  $2^K K!$  equivalent modes in the posterior, a subset of which is explored by the sampler. This causes most posterior summaries (e.g. posterior mean) to be of little utility.

A relabelling algorithm, similar to that of Erosheva and Curtis (2017), is used to deal with these issues of label switching and sign switching. Following the method of Stephens (2000), a decision-theoretic approach is to define a loss function for a set of actions and relabellings, and select the action and relabelling which minimises the posterior expected loss. This is done with the aim of relabelling samples such that they correspond to being sampled around the same mode.

Define an action

$$\mathbf{a} = \left( \{m_{l_{ik}}\}_{i=1:G, k=1:K}, \{s_{l_{ik}}^2\}_{i=1:G, k=1:K}, \{m_{f_{kj}}\}_{j=1:N, k=1:K}, \{s_{f_{kj}}^2\}_{j=1:N, k=1:K}, \{p_{z_{ik}}\}_{i=1:G, k=1:K} \right)$$

to be a choice of means and variances of the entries of  $\mathbf{L}$  and  $\mathbf{F}$ , and also the means of the entries of  $\mathbf{Z}$ . Let  $\sigma \in S_K$  and  $\boldsymbol{\nu} \in \{-1, 1\}^K$ , where  $S_K$  is the set of permutations on the set  $\{1, 2, \dots, K\}$ . We define a loss function as follows:

$$\begin{aligned} \mathcal{L}(\mathbf{a}, \sigma, \boldsymbol{\nu}; \mathbf{L}, \mathbf{F}, \mathbf{Z}) = & - \sum_{k=1}^K \left\{ \sum_{i=1}^G \log \text{Bern}(z_{i\sigma(k)} \mid p_{z_{ik}}) + \mathbb{I}(l_{i\sigma(k)} \neq 0) \log \mathcal{N}(\nu_k l_{i\sigma(k)} \mid m_{l_{ik}}, s_{l_{ik}}^2) \right. \\ & \left. + \sum_{j=1}^N \mathcal{N}(\nu_k f_{\sigma(k)j} \mid m_{f_{kj}}, s_{f_{kj}}^2) \right\}. \end{aligned}$$

Suppose we want to relabel  $T$  samples  $\{(\mathbf{L}^{(t)}, \mathbf{F}^{(t)}, \mathbf{Z}^{(t)})\}_{t=1:T}$  obtained from MCMC. We seek to choose  $\mathbf{a}$  and  $\{(\sigma^{(t)}, \boldsymbol{\nu}^{(t)})\}_{t=1:T}$  such that the Monte Carlo risk

$$\mathcal{R}_{\text{MC}} = \sum_{t=1}^T \mathcal{L}(\mathbf{a}, \{(\sigma^{(t)}, \boldsymbol{\nu}^{(t)})\}_{t=1:T}; \mathbf{L}^{(t)}, \mathbf{F}^{(t)}, \mathbf{Z}^{(t)})$$

is minimised. After initialising  $\mathbf{a}$  and  $\{(\sigma^{(t)}, \boldsymbol{\nu}^{(t)})\}_{t=1:T}$ , a local optimum may be obtained by alternating between the following steps:

1. Choose  $\mathbf{a}$  such that the Monte Carlo risk is minimised given the current values of  $\{(\sigma^{(t)}, \boldsymbol{\nu}^{(t)})\}_{t=1:T}$ .
2. Choose  $\{(\sigma^{(t)}, \boldsymbol{\nu}^{(t)})\}_{t=1:T}$  such that the Monte Carlo risk is minimised given the current action  $\mathbf{a}$ .

The procedure is terminated when a fixed point is reached.

Step 1 may be solved analytically, by setting partial derivatives of the Monte Carlo risk with respect to the action parameters to zero. This is equivalent to finding the maximum likelihood estimators, summarised by the following updates:

$$\begin{aligned}
\widehat{m}_{l_{ik}} &= \frac{\sum_{t=1}^T \mathbb{I}(l_{i\sigma^{(t)}(k)} \neq 0) \nu_k^{(t)} l_{i\sigma^{(t)}(k)}^{(t)}}{\sum_{t=1}^T \mathbb{I}(l_{i\sigma^{(t)}(k)} \neq 0)} \\
\widehat{s}_{l_{ik}}^2 &= \frac{\sum_{t=1}^T \mathbb{I}(l_{i\sigma^{(t)}(k)} \neq 0) \left( \nu_k^{(t)} l_{i\sigma^{(t)}(k)}^{(t)} - \widehat{m}_{l_{ik}} \right)^2}{\sum_{t=1}^T \mathbb{I}(l_{i\sigma^{(t)}(k)} \neq 0)} \\
\widehat{m}_{f_{kj}} &= \frac{\sum_{t=1}^T \mathbb{I}(f_{\sigma^{(t)}(k)j} \neq 0) \nu_k^{(t)} f_{\sigma^{(t)}(k)j}^{(t)}}{\sum_{t=1}^T \mathbb{I}(f_{\sigma^{(t)}(k)j} \neq 0)} \\
\widehat{s}_{f_{kj}}^2 &= \frac{\sum_{t=1}^T \mathbb{I}(f_{\sigma^{(t)}(k)j} \neq 0) \left( \nu_k^{(t)} f_{\sigma^{(t)}(k)j}^{(t)} - \widehat{m}_{f_{kj}} \right)^2}{\sum_{t=1}^T \mathbb{I}(f_{\sigma^{(t)}(k)j} \neq 0)} \\
\widehat{p}_{z_{ik}} &= \frac{\sum_{t=1}^T z_{i\sigma^{(t)}(k)}^{(t)}}{T}.
\end{aligned}$$

Step 2 is equivalent to the linear assignment problem, which may be solved by a  $\mathcal{O}(K^3)$  algorithm of Jonker and Volgenant (1987).

### 3 Variational Inference

#### 3.1 Mean-field approximation

Use the variational factorisation

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

as an approximation to the posterior distribution, where

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

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}, \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}, \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{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'}} \mu_{l_{ik'}} l_{ik} \right. \right. \\
&\quad \left. \left. + z_{ik} \sum_{j=1}^N \left( \left[ \Sigma_{\mathbf{f}_{\cdot j}} \right]_{kk} + \left[ \boldsymbol{\mu}_{\mathbf{f}_{\cdot j}} \right]_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

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

$$\begin{aligned}
\mu_{l_{ik}}^* &= \frac{\hat{a}_{\tau_i}}{\hat{b}_{\tau_i}} \sigma_{l_{ik}}^{2*} \boldsymbol{\mu}_{\mathbf{f}_k}^\top \left( \mathbf{y}_i - \sum_{k' \neq k} \eta_{ik'} \boldsymbol{\mu}_{\mathbf{f}_{k'}} \mu_{l_{ik'}} \right) \\
q(z_{ik}) &\propto \exp \left\{ \frac{z_{ik}}{2} \left( \psi(\hat{a}_{\alpha_k}) - \log 2\pi \hat{b}_{\alpha_k} + \frac{\mu_{l_{ik}}^{2*}}{\sigma_{l_{ik}}^{2*}} \right) \right\} \left( \sqrt{\sigma_{l_{ik}}^{2*} \pi_k} \right)^{z_{ik}} (1 - \pi_k)^{1-z_{ik}}. \quad (9)
\end{aligned}$$

Coordinate ascent for  $\mathbf{f}_{\cdot j}$  gives

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

where

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

which corresponds to the updates

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

Coordinate ascent for  $\tau_i$  gives

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

where

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

which corresponds to the updates

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

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

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

## References

- Erosheva, E. A. and Curtis, S. M. (2017). Dealing with Reflection Invariance in Bayesian Factor Analysis. *Psychometrika*, 82(2):295–307.
- Jonker, R. and Volgenant, A. (1987). A shortest augmenting path algorithm for dense and sparse linear assignment problems. *Computing*, 38(4):325–340.
- Stephens, M. (2000). Dealing with label switching in mixture models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 62(4):795–809.

## A Unused attempts

### A.1 Capturing dependency within $\mathbf{l}_i$ .

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

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

where

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

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, \boldsymbol{\alpha}} [\log p(\mathbf{l}_i, \mathbf{z}_i \mid \mathbf{Y}, \mathbf{F}, \boldsymbol{\tau}, \boldsymbol{\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. \left. - \frac{1}{2} \mathbb{E}_{\boldsymbol{\alpha}} \left[ [\mathbf{l}_i]_{\mathbf{z}_i}^\top [D\boldsymbol{\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. \left. - \frac{1}{2} [\mathbf{l}_i]_{\mathbf{z}_i}^\top [D\boldsymbol{\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\boldsymbol{\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

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

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

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

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

Marginalising out  $\mathbf{l}_{i\cdot}$  then gives

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

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

$$\hat{q}(\mathbf{z}_{i\cdot}) = \prod_{k=1}^K \hat{q}(z_{ik}) = \prod_{k=1}^K \text{Bern}(z_{ik} \mid \gamma_{ik}).\tag{16}$$

Methods for estimating  $\gamma_{ik}$  will be addressed in Section A.2.

For the remaining variational parameters, computations similar to the Section 3.1 give the following updates:

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

$$\boldsymbol{\mu}_{\mathbf{f}\cdot j}^* = \Sigma_{\mathbf{f}\cdot j}^* \overline{\mathbf{L}}^\top D_{\boldsymbol{\tau}} \mathbf{y}_{\cdot j}\tag{18}$$

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

$$\hat{b}_{\tau_i}^* = b_{\tau} + \frac{1}{2} \left( \mathbf{y}_{i\cdot}^\top \mathbf{y}_{i\cdot} - 2 \overline{\mathbf{l}_{i\cdot}}^\top \overline{\mathbf{F}} \mathbf{y}_{i\cdot} + \overline{\mathbf{l}_{i\cdot}}^\top \overline{\mathbf{F} \mathbf{F}^\top} \mathbf{l}_{i\cdot} \right)\tag{18}$$

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

$$\hat{b}_{\alpha_k}^* = b_{\alpha} + \frac{1}{2} \sum_{i=1}^G \eta_{ik} (\sigma_{l_{ik}}^2 + \mu_{l_{ik}}^2).\tag{19}$$

However, these expectations depend on  $\boldsymbol{\mu}_{\mathbf{l}_{i\cdot}}$  and  $\Sigma_{\mathbf{l}_{i\cdot}}$ , which we do not have values of due to their varying dimensions.

## A.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

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



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  $q^*(\zeta)$  for  $\zeta \in \mathbf{B}_T$ , where  $\mathbf{B}_T$  is a random subset of  $B$  of size  $T$ .

Define

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

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:

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

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.