

# Infinite Mixtures of Infinite Factor Analysers

## Notes & Derivations

Keefe Murphy<sup>1, 2</sup>, Dr. Claire Gormley<sup>1, 2</sup>, and Prof. Brendan Murphy<sup>1, 2</sup>

<sup>1</sup>Department of Mathematics and Statistics, UCD

<sup>2</sup>Insight Centre for Data Analytics, UCD

# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>                                     | <b>3</b>  |
| 1.1      | Model Set-Up . . . . .                                  | 3         |
| 1.2      | Assumptions . . . . .                                   | 3         |
| <b>2</b> | <b>Bayesian Framework</b>                               | <b>4</b>  |
| 2.1      | Likelihood . . . . .                                    | 4         |
| 2.2      | Posterior Set-Up . . . . .                              | 5         |
| <b>3</b> | <b>Sampling from the Full Conditionals</b>              | <b>5</b>  |
| 3.1      | Factor Scores . . . . .                                 | 5         |
| 3.2      | Loadings Matrix . . . . .                               | 6         |
| 3.3      | Uniquenesses . . . . .                                  | 7         |
| 3.4      | Reintroducing $\mu$ . . . . .                           | 7         |
| 3.5      | Gibbs Sampler Pseudo-Code . . . . .                     | 8         |
| 3.6      | Issues Around Identifiability . . . . .                 | 9         |
| <b>4</b> | <b>Introducing the Shrinkage Prior</b>                  | <b>9</b>  |
| 4.1      | Multiplicative Gamma Process Shrinkage Priors . . . . . | 9         |
| 4.2      | Deriving new MGP Full Conditionals . . . . .            | 10        |
| 4.2.1    | Loadings Matrix . . . . .                               | 10        |
| 4.2.2    | Local Shrinkage . . . . .                               | 10        |
| 4.2.3    | Global Shrinkage . . . . .                              | 11        |
| 4.3      | Adaptive Step . . . . .                                 | 11        |
| <b>5</b> | <b>Extension to Clustering Heterogeneous Data</b>       | <b>12</b> |
| 5.1      | Introducing Mixture Models . . . . .                    | 12        |
| 5.1.1    | Decomposable Prior for $\gamma$ . . . . .               | 13        |
| 5.2      | Deriving Posterior Distributions . . . . .              | 13        |
| 5.2.1    | Cluster Mixing Proportions . . . . .                    | 14        |
| 5.2.2    | Cluster Labels . . . . .                                | 14        |
| 5.2.3    | MCMC Algorithm Pseudo-Code . . . . .                    | 14        |
| 5.3      | Label Switching . . . . .                               | 15        |
| <b>6</b> | <b>References</b>                                       | <b>15</b> |

# 1 Introduction

## 1.1 Model Set-Up

Let  $\underline{x} = (x_1, x_2, \dots, x_p)^T$  have mean  $\underline{\mu}$  and covariance matrix  $\underline{\Sigma}$ . The factor model states that  $\underline{x}$  is linearly independent upon a few ( $q \ll p$ ) unobservable random variables  $f_1, f_2, \dots, f_q$  called *common factors* and  $p$  additional sources of variation  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_p$  called *specific factors*, s.t.

$$\underline{x}_i = \underline{\mu} + \underline{\Lambda} \underline{f}_i + \underline{\varepsilon}_i$$

|       |                             |               |                |                                     |
|-------|-----------------------------|---------------|----------------|-------------------------------------|
| where | $\underline{x}_i$           | $\rightarrow$ | $(p \times 1)$ | observation vector                  |
|       | $\underline{\mu}$           | $\rightarrow$ | $(p \times 1)$ | overall mean vector                 |
|       | $\underline{\Lambda}$       | $\rightarrow$ | $(p \times q)$ | loadings matrix                     |
|       | $\underline{f}_i$           | $\rightarrow$ | $(q \times 1)$ | vector of factor scores for obs $i$ |
|       | $\underline{\varepsilon}_i$ | $\rightarrow$ | $(p \times 1)$ | vector of errors for obs $i$        |
|       | $i$                         | $=$           | $1, \dots, N$  |                                     |
|       | $j$                         | $=$           | $1, \dots, p$  |                                     |
|       | $k$                         | $=$           | $1, \dots, q$  |                                     |

$\Lambda_{jk}$  is called the *factor loading* of the  $j$ -th variable on the  $k$ -th factor of the  $(p \times q)$   $\underline{\Lambda}$  matrix of factor loadings. If we assume the data has been centred to have column means of 0 then we have

$$\left( \underline{x}_i - \underline{\mu} \right)_{(p \times 1)} = \underline{x}_{i(p \times 1)}^* = \underline{\Lambda}_{(p \times q)} \underline{f}_{i(q \times 1)} + \underline{\varepsilon}_{i(p \times 1)} \quad (1)$$

## 1.2 Assumptions

1.  $\underline{\mu} = 0$
2.  $\underline{\varepsilon}_i$  and  $\underline{f}_i$  are independent:  $\text{Cov}(\underline{f}, \underline{\varepsilon}) = \text{E}(\underline{f}, \underline{\varepsilon}^T) = 0$
3.  $\underline{\varepsilon}_i \sim \text{N}(\underline{0}, \underline{\Psi})$  where  $\underline{\Psi} = \text{diag}(\psi_1, \dots, \psi_p)$

$$\therefore \text{E}(\underline{\varepsilon}) = \underline{0} \text{ and } \text{Cov}(\underline{\varepsilon}) = \begin{pmatrix} \psi_1 & 0 & \dots & 0 \\ 0 & \psi_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \psi_p \end{pmatrix} = \underline{\Psi}$$

$$\therefore \underline{\varepsilon}_i \sim \text{MVN}_p(\underline{0}, \underline{\Psi}) \quad (2)$$

4.  $\underline{f}_i \sim \text{MVN}_q(\underline{0}, \underline{\mathcal{I}}_q)$

$$\therefore \text{E}(\underline{f}) = \underline{0} \text{ and } \text{Cov}(\underline{f}) = \begin{pmatrix} 1_1 & 0 & \dots & 0 \\ 0 & 1_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1_q \end{pmatrix} = \underline{\mathcal{I}}_q$$

$$\therefore \underline{f}_i \sim \text{MVN}_q(\underline{0}, \underline{\mathcal{I}}_q) \quad (3)$$

## 2 Bayesian Framework

### 2.1 Likelihood

$$\begin{aligned}
\mathbb{E}(\underline{x}_i^*) &= \mathbb{E}(\underline{\Lambda}\underline{f}_i + \underline{\varepsilon}_i) \\
&= \underline{\Lambda}\mathbb{E}(\underline{f}_i) + \mathbb{E}(\underline{\varepsilon}_i) \\
&= \underline{0} \\
\therefore \underline{X}_i^* &\sim \text{MVN}_p(\underline{0}, \underline{\Sigma})
\end{aligned} \tag{4}$$

$$\begin{aligned}
\text{Since } \underline{\varepsilon}_i &= \underline{x}_i^* - \underline{\Lambda}\underline{f}_i \\
\underline{\Sigma} &= \text{Cov}(\underline{X}) \\
&= \mathbb{E}\left[(\underline{x} - \underline{\mu})(\underline{x} - \underline{\mu})^T\right] \\
&= \mathbb{E}\left[\underline{x}^* \underline{x}^{*T}\right] \\
&= \mathbb{E}\left[(\underline{\Lambda}\underline{f} + \underline{\varepsilon})(\underline{\Lambda}\underline{f} + \underline{\varepsilon})^T\right] \\
&= \mathbb{E}\left[(\underline{\Lambda}\underline{f}) + \underline{\varepsilon}(\underline{\Lambda}\underline{f})^T + (\underline{\Lambda}\underline{f})\underline{\varepsilon}^T + \underline{\varepsilon}\underline{\varepsilon}^T\right] \\
&= \underline{\Lambda}\mathbb{E}(\underline{f}\underline{f}^T)\underline{\Lambda}^T + \mathbb{E}(\underline{\varepsilon}\underline{f}^T)\underline{\Lambda}^T + \underline{\Lambda}\mathbb{E}(\underline{f}\underline{\varepsilon}^T) + \mathbb{E}(\underline{\varepsilon}\underline{\varepsilon}^T) \\
&= \underline{\Lambda}\underline{\Lambda}^T + \underline{\Psi} \\
\therefore \underline{X}_i^* &\sim \text{MVN}_p(\underline{0}, \underline{\Lambda}\underline{\Lambda}^T + \underline{\Psi})
\end{aligned} \tag{5}$$

$$\begin{aligned}
\mathbb{E}(\underline{X}_i^* | \underline{f}_i) &= \mathbb{E}(\underline{\Lambda}\underline{f}_i + \underline{\varepsilon}_i | \underline{f}_i) \\
&= \underline{\Lambda}\mathbb{E}(\underline{f}_i | \underline{f}_i) + \mathbb{E}(\underline{\varepsilon}_i | \underline{f}_i) \\
&= \underline{\Lambda}\underline{f}_i \\
\text{Cov}(\underline{X}_i^* | \underline{f}_i) &= \mathbb{E}\left[(\underline{X}_i^* - \underline{\Lambda}\underline{f}_i)(\underline{X}_i^* - \underline{\Lambda}\underline{f}_i)^T | \underline{f}_i\right] \\
&= \mathbb{E}(\underline{\varepsilon}_i \underline{\varepsilon}_i^T | \underline{f}_i) \\
&= \underline{\Psi} \\
\therefore \underline{X}_i^* | \underline{f}_i, \underline{\Lambda}, \underline{\Psi} &\sim \text{MVN}_p(\underline{\Lambda}\underline{f}_i, \underline{\Psi})
\end{aligned} \tag{6}$$

The density of the data is then given by:

$$\begin{aligned}
P(\underline{X}_i^* | \underline{f}_i, \underline{\Lambda}, \underline{\Psi}) &= (2\pi)^{-\frac{p}{2}} |\underline{\Psi}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \sum_{i=1}^N (\underline{X}_i^* - \underline{\Lambda}\underline{f}_i)^T \underline{\Psi}^{-1} (\underline{X}_i^* - \underline{\Lambda}\underline{f}_i)\right) \\
&\propto |\underline{\Psi}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \text{tr}\left[\underline{\Psi}^{-1} (\underline{X} - \underline{F}\underline{\Lambda})^T (\underline{X} - \underline{F}\underline{\Lambda})\right]\right)
\end{aligned} \tag{7}$$

$$\begin{aligned}
\text{Where } \underline{\Lambda}_{(p \times q)} &= \begin{pmatrix} \lambda_{11} & \lambda_{12} & \dots & \lambda_{1q} \\ \lambda_{21} & \lambda_{22} & \dots & \lambda_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{p1} & \lambda_{p2} & \dots & \lambda_{pq} \end{pmatrix} \\
&\& \underline{F}_{(n \times q)} = \begin{pmatrix} f_{11} & f_{12} & \dots & f_{1q} \\ f_{21} & f_{22} & \dots & f_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ f_{n1} & f_{n2} & \dots & f_{nq} \end{pmatrix} \& \underline{f}_i \text{ is a column vector containing the entries of row } i \text{ of } \underline{F}
\end{aligned}$$

## 2.2 Posterior Set-Up

$$\begin{aligned}
\text{Likelihood} &= P(X^* | \underline{\theta}) \\
&= P(\underline{X}_i^* | \underline{f}_i, \underline{\Lambda}, \underline{\Psi}) \\
\therefore P(\underline{X}_i^* | \underline{f}_i, \underline{\Lambda}, \underline{\Psi}) &\sim \text{MVN}_p(\underline{\Lambda} \underline{f}_i, \underline{\Psi})
\end{aligned} \tag{8}$$

$$\begin{aligned}
\text{Prior} &= P(\underline{\theta}) \\
&= P(\underline{F}) P(\underline{\Lambda}) P(\underline{\Psi})
\end{aligned}$$

$$\begin{aligned}
\text{Posterior} &\propto \text{Likelihood} \times \text{Prior} \\
\therefore P(\underline{F}, \underline{\Lambda}, \underline{\Psi} | \underline{X}) &\propto \mathcal{L}(\underline{X}^* | \underline{F}, \underline{\Lambda}, \underline{\Psi}) P(\underline{F}) P(\underline{\Lambda}) P(\underline{\Psi}) \\
&\propto \left[ \prod_{i=1}^N P(\underline{X}_i^* | \underline{f}_i, \underline{\Lambda}, \underline{\Psi}) \right] \left[ \prod_{i=1}^N P(\underline{f}_i) \right] \left[ \prod_{j=1}^p P(\underline{\Lambda}_j) \right] \left[ \prod_{j=1}^p P(\Psi_{jj}) \right]
\end{aligned} \tag{9}$$

Later on, especially as we move into the mixture case, it will be necessary to undo the centering, thereby removing the  $*$  on  $\underline{X}$ , and reintroduce  $\underline{\mu}$ . This will necessitate multiplying the quantity in (Equation (9)) by  $\left[ \prod_{j=1}^p P(\underline{\mu}_j) \right]$ . However, we will proceed to derive the full conditionals we need for Gibbs Sampling using the centered notation for now as adjusting for  $\underline{\mu}$  afterwards will be trivial.

## 3 Sampling from the Full Conditionals

### 3.1 Factor Scores - $\underline{f}_i$

$$\begin{aligned}
\underline{f}_i &\sim \text{MVN}_q(0, \mathcal{I}_q) \\
&= (2\pi)^{-\frac{q}{2}} \exp\left(-\frac{1}{2} \underline{f}_i^T \underline{f}_i\right)
\end{aligned} \tag{10}$$

To obtain the full conditional for  $\underline{f}_i$  we can multiply the conditional likelihood by the marginal distribution in (Equation(10)) s.t.

$$\begin{aligned}
P(\underline{f}_i | \underline{X}_i^*, \underline{\Lambda}, \underline{\Psi}) &\sim P(\underline{X}_i^* | \underline{f}_i, \underline{\Lambda}, \underline{\Psi}) P(\underline{f}_i) \\
&\propto \exp\left(-\frac{1}{2} \sum_{i=1}^N (\underline{X}_i^* - \underline{\Lambda} \underline{f}_i)^T \underline{\Psi}^{-1} (\underline{X}_i^* - \underline{\Lambda} \underline{f}_i) + \underline{f}_i^T \underline{f}_i\right) \\
&\propto \exp\left(-\frac{1}{2} \sum_{i=1}^N \left[-\underline{X}_i^{*T} \underline{\Psi}^{-1} \underline{\Lambda} \underline{f}_i - (\underline{\Lambda} \underline{f}_i)^T \underline{\Psi}^{-1} \underline{X}_i^* + (\underline{\Lambda} \underline{f}_i)^T \underline{\Psi}^{-1} (\underline{\Lambda} \underline{f}_i) - \underline{f}_i^T \underline{f}_i\right]\right) \\
&\propto \exp\left(-\frac{1}{2} \left\{ \underline{f}_i^T [\mathcal{I}_q + \underline{\Lambda}^T \underline{\Psi}^{-1} \underline{\Lambda}] \underline{f}_i \right\} + \underline{X}_i^{*T} \underline{\Psi}^{-1} \underline{\Lambda} \underline{f}_i\right)
\end{aligned} \tag{11}$$

As this is the product of two MVN distributions we can expect the result to also be MVN. Typically,

$$\begin{aligned}
\text{MVN}(x; \underline{\mu}, \underline{\Sigma}) &\propto \exp\left(-\frac{1}{2} (\underline{X} - \underline{\mu})^T \underline{\Sigma}^{-1} (\underline{X} - \underline{\mu})\right) \\
&= \exp\left(-\frac{1}{2} (\underline{X}^T \underline{\Sigma}^{-1} \underline{X} - 2 \underline{\mu}^T \underline{\Sigma}^{-1} \underline{X} + \underline{\mu}^T \underline{\Sigma}^{-1} \underline{\mu})\right)
\end{aligned}$$

$\therefore$  we can identify the  $\underline{\mu}$  and  $\underline{\Sigma}^{-1}$  terms from (Equation (11)) above to yield

$$P(\underline{f}_i | \underline{X}_i^*, \underline{\Lambda}, \underline{\Psi}) \sim \text{MVN}_q\left([\mathcal{I}_q + \underline{\Lambda}^T \underline{\Psi}^{-1} \underline{\Lambda}]^{-1} \underline{\Lambda}^T \underline{\Psi}^{-1} \underline{X}_i^*, [\mathcal{I}_q + \underline{\Lambda}^T \underline{\Psi}^{-1} \underline{\Lambda}]^{-1}\right) \tag{12}$$

However, we can reintroduce  $\underline{\mu}$  and save on computational time if we implement the algorithm of Rue & Held (2005)<sup>1</sup>. In fact, we can extend this to block update the scores, thereby obviating the need to loop over  $i$ :

- Calculate  $\underline{\Omega}_F = \mathcal{I}_q + \underline{\Lambda}^T \underline{\Psi}^{-1} \underline{\Lambda}$
- Compute the Cholesky Factorization  $\underline{\Omega}_F = U^T U$ .
- Sample  $z \sim \text{MVN}_q(0, \mathcal{I}_q)$   $N$  times.
- Backsolve  $Uv = z^T$ .
- Compute  $\underline{\Omega}_F^{-1}$  from  $U$ .
- Return  $\left( \underline{\Omega}_F^{-1} \underline{\Lambda}^T \underline{\Psi}^{-1} (\underline{X} - \underline{\mu})^T + v \right)^T$  (13).

### 3.2 Loadings Matrix - $\underline{\Lambda}$

A Gaussian distribution is a conjugate prior for  $\underline{\Lambda}$ , implying an  $\text{MVN}_q$  distribution prior for each row  $\underline{\Lambda}_j$  of  $\underline{\Lambda}$  s.t.  $\underline{\Lambda}_j \sim \text{MVN}_q(\underline{0}, \Sigma_\lambda \mathcal{I}_q)$  where  $\Sigma_\lambda$  is a scalar hyperparameter which controls the diagonal covariance matrix of the prior. As above, we can expect the result of the product of two  $\text{MVN}_q$  distributions to itself be distributed in the same way.

$$\begin{aligned}
P(\underline{\Lambda}_j | \underline{X}^*, \underline{F}, \underline{\Psi}) &\sim P(\underline{X}^* | \underline{F}, \underline{\Lambda}_j, \underline{\Psi}) P(\underline{\Lambda}_j | \Sigma_\lambda) \\
&\propto \exp \left( -\frac{1}{2} \sum_{i=1}^N (\underline{X}_i^* - \underline{\Lambda}_j \underline{f}_i)^T \Psi_{jj}^{-1} (\underline{X}_i^* - \underline{\Lambda}_j \underline{f}_i) \right) \exp \left( -\frac{1}{2} (\underline{\Lambda}_j^T (\Sigma_\lambda \mathcal{I}_q)^{-1} \underline{\Lambda}_j) \right) \\
&\propto \exp \left( -\frac{1}{2} \sum_{i=1}^N \left[ -2 \underline{X}_i^{*T} \Psi_{jj}^{-1} (\underline{\Lambda}_j \underline{f}_i) + (\underline{\Lambda}_j \underline{f}_i)^T \Psi_{jj}^{-1} (\underline{\Lambda}_j \underline{f}_i) + \underline{\Lambda}_j^T (\Sigma_\lambda \mathcal{I}_q)^{-1} \underline{\Lambda}_j \right] \right) \\
&\propto \exp \left( \underline{\Lambda}_j \Psi_{jj}^{-1} \sum_{i=1}^N \underline{X}_{ij}^{*T} \underline{f}_i - \frac{1}{2} \underline{\Lambda}_j^T \left[ \sum_{i=1}^N \Psi_{jj}^{-1} \underline{f}_i^T \underline{f}_i \right] \underline{\Lambda}_j - \frac{1}{2} \underline{\Lambda}_j^T (\Sigma_\lambda \mathcal{I}_q)^{-1} \underline{\Lambda}_j \right) \\
&\propto \exp \left( \underline{\Lambda}_j [\underline{F}^T \Psi_{jj}^{-1} \underline{X}^{j*}] - \frac{1}{2} \underline{\Lambda}_j^T [(\Sigma_\lambda \mathcal{I}_q)^{-1} + \Psi_{jj}^{-1} \underline{F}^T \underline{F}] \underline{\Lambda}_j \right) \quad (14)
\end{aligned}$$

where  $\underline{X}^{j*}$  denotes the  $j$ -th column of  $\underline{X}^*$

$$\therefore P(\underline{\Lambda}_j | \underline{X}^*, \underline{F}, \underline{\Psi}) \sim \text{MVN}_q \left( [(\Sigma_\lambda \mathcal{I}_q)^{-1} + \Psi_{jj}^{-1} \underline{F}^T \underline{F}]^{-1} \underline{F}^T \Psi_{jj}^{-1} \underline{X}^{j*}, [(\Sigma_\lambda \mathcal{I}_q)^{-1} + \Psi_{jj}^{-1} \underline{F}^T \underline{F}]^{-1} \right) \quad (15)$$

However, we can reintroduce  $\underline{\mu}$  and save on computational time, as before, if we:

- Calculate  $\underline{\Omega}_{\lambda_j} = (\Sigma_\lambda \mathcal{I}_q)^{-1} + \Psi_{jj}^{-1} \underline{F}^T \underline{F}$ .
- Compute the Cholesky Factorization  $\underline{\Omega}_{\lambda_j} = U^T U$ .
- Sample  $z \sim \text{N}(0, 1)$   $q$  times.
- Back-solve  $Uv = z$ .
- Compute  $\underline{\Omega}_{\lambda_j}^{-1}$  from  $U$ .
- Return  $\underline{\Omega}_{\lambda_j}^{-1} \underline{F}^T \Psi_{jj}^{-1} (\underline{X}^j - \underline{\mu}_j) + v$  (16).

<sup>1</sup>To sample  $x \sim \text{N}(\mu, \Omega^{-1})$ , find a matrix  $U$  – non-unique, and square or ‘tall’ – via Cholesky Decomposition s.t.  $U^T U = \Omega$ , sample from  $z \sim \text{N}(0, 1)$ , then backsolve  $L^T v = Uv = z$  s.t.  $x = \mu + v = \mu + L^{-T} z = \mu + U^{-1} z$ . Then:

- $E(x) = \mu + U^{-1} E(z) = \mu$
- $\text{Cov}(x, x) = \text{Cov}(L^{-T} z, z) = (L^T L)^{-1} = \Omega^{-1}$

### 3.3 Uniquenesses - $\underline{\Psi}$

If we suggest an Inverse Wishart prior distribution for  $\underline{\Psi}$ , we have:

$$P(\underline{\Psi}) \propto |\underline{\Psi}^{-1}|^{\frac{N+p+1}{2}} \exp\left(-\frac{1}{2}\text{tr}(\underline{\mathcal{S}}^{-1*}\underline{\Psi})\right)$$

Using the fact that  $V^{-1} \sim \text{Wish}_p(\nu, \Sigma)$  when  $V \sim \text{Wish}_p^{-1}(m, \Sigma^{-1})$  with  $m = \nu + p + 1$  we can rewrite as:

$$P(\underline{\Psi}^{-1}) \propto |\underline{\Psi}^{-1}|^{\frac{N}{2}} \exp\left(-\frac{1}{2}\text{tr}(\underline{\mathcal{S}}^*\underline{\Psi}^{-1})\right)$$

Since  $\underline{\Psi}$  is a diagonal matrix:

$$P(\underline{\Psi}^{-1}) \propto \prod_{j=1}^p |\Psi_{jj}^{-1}|^{\frac{N}{2}} \exp\left(-\frac{1}{2}\text{tr}(\mathcal{S}_{jj}^*\Psi_{jj}^{-1})\right)$$

This suggests the prior for  $\underline{\Psi}^{-1}$  is a product of  $p$   $\text{Ga}(\alpha/2, \beta/2)$  distributions.

$$\begin{aligned} \therefore P(\underline{\Psi}^{-1} | \alpha, \beta) &= \prod_{j=1}^p P(\Psi_{jj}^{-1} | \alpha, \beta) \\ &\propto \prod_{j=1}^p (\Psi_{jj}^{-1})^{\frac{\alpha}{2}-1} \exp\left(-\frac{\beta}{2}\Psi_{jj}^{-1}\right) \\ \therefore P(\underline{\Psi}^{-1} | \underline{X}^*, \underline{F}, \underline{\Lambda}) &\propto P(\underline{X}^* | \underline{F}, \underline{\Lambda}) P(\underline{\Psi}^{-1} | \alpha, \beta) \\ &\propto \prod_{j=1}^p (\Psi_{jj}^{-1})^{\frac{N}{2}} \exp\left(-\frac{\mathcal{S}_{jj}^*}{2}\Psi_{jj}^{-1}\right) \prod_{j=1}^p (\Psi_{jj}^{-1})^{\frac{\alpha}{2}-1} \exp\left(-\frac{\beta}{2}\Psi_{jj}^{-1}\right) \\ &\propto \prod_{j=1}^p (\Psi_{jj}^{-1})^{\frac{N+\alpha}{2}-1} \exp\left(-\frac{\mathcal{S}_{jj}^* + \beta}{2}\Psi_{jj}^{-1}\right) \end{aligned} \quad (17)$$

$$\text{where } \mathcal{S}_{jj}^* = \sum_{i=1}^N (x_{ij} - \underline{\Lambda}_j \mathbf{f}_i)^2$$

However, we can reintroduce  $\underline{\mu}$  at this stage by rewriting:

$$\mathcal{S}_{jj} = \sum_{i=1}^N (x_{ij} - \mu_j - \underline{\Lambda}_j \mathbf{f}_i)^2$$

Thus the posterior distribution of each  $\Psi_{jj}^{-1}$  is given by:

$$P(\Psi_{jj}^{-1} | \underline{X}, \underline{F}, \underline{\Lambda}) \sim \text{Ga}\left(\frac{N + \alpha}{2}, \frac{\mathcal{S}_{jj} + \beta}{2}\right) \quad (18)$$

### 3.4 Reintroducing $\underline{\mu}$

We've already seen from (Equations (13), (16) and (18)) that reintroducing  $\underline{\mu}$  to the other full conditionals is trivial. All that remains is to specify the conjugate Gaussian prior for  $\underline{\mu}$  itself, and to derive its full conditional. This implies an  $\text{MVN}_p$  distribution prior s.t.  $\underline{\mu} \sim \text{MVN}_p(\underline{0}, \Sigma_{\mu} \mathcal{I}_p)$  where  $\Sigma_{\mu}$  is a scalar hyperparameter which controls the diagonal covariance matrix of the prior. As above, we can expect the result of the product of two  $\text{MVN}_p$  distributions to itself be distributed in the same way.

$$\begin{aligned} P(\underline{\mu} | \underline{X}, \underline{F}, \underline{\Psi}, \underline{\Lambda}) &\propto \exp\left(-\frac{1}{2} \sum_{i=1}^N (\underline{X}_i - \underline{\mu} - \underline{\Lambda} \mathbf{f}_i)^T \underline{\Psi}^{-1} (\underline{X}_i - \underline{\mu} - \underline{\Lambda} \mathbf{f}_i)\right) \exp\left(-\frac{1}{2} (\underline{\mu}^T (\Sigma_{\mu} \mathcal{I}_p)^{-1} \underline{\mu})\right) \\ &\propto \exp\left(-\frac{1}{2} \sum_{i=1}^N \left[-2 \underline{X}_i^T \underline{\Psi}^{-1} \underline{\mu} + 2 (\underline{\Lambda} \mathbf{f}_i)^T \underline{\Psi}^{-1} \underline{\mu} + \underline{\mu}^T \underline{\Psi}^{-1} \underline{\mu} + \underline{\mu}^T (\Sigma_{\mu} \mathcal{I}_p)^{-1} \underline{\mu}\right]\right) \\ &\propto \exp\left(\sum_{i=1}^N \underline{X}_i^T \underline{\Psi}^{-1} \underline{\mu} - \sum_{i=1}^N (\underline{\Lambda} \mathbf{f}_i)^T \underline{\Psi}^{-1} \underline{\mu} - \frac{1}{2} [\underline{\mu}^T ((\Sigma_{\mu} \mathcal{I}_p)^{-1} + N \underline{\Psi}^{-1}) \underline{\mu}]\right) \end{aligned} \quad (19)$$

$$\therefore P(\underline{\mu} | \underline{X}, \underline{F}, \underline{\Psi}, \underline{\Lambda}) \sim \text{MVN}_p \left( \left[ (\Sigma_\mu \mathcal{I}_p)^{-1} + n \underline{\Psi}^{-1} \right]^{-1} \underline{\Psi}^{-1} \left( \sum_{i=1}^N \underline{X}_i - \sum_{i=1}^N \underline{\Lambda} \underline{f}_i \right), \right. \\ \left. \left[ (\Sigma_\mu \mathcal{I}_p)^{-1} + N \underline{\Psi}^{-1} \right]^{-1} \right) \quad (20)$$

However, we can save on computational time, as before, if we:

- Calculate  $\underline{\Omega}_\mu = (\Sigma_\mu \mathcal{I}_p)^{-1} + N \underline{\Psi}^{-1}$ , which is a diagonal  $p \times p$  matrix.
- Invert  $\underline{\Omega}_\mu$  by inverting its diagonal elements.
- Compute the Cholesky Factorization  $\underline{\Omega}_\mu^{-1} = U^T U$ .
- Sample  $z \sim N(0, 1)$   $q$  times.
- Compute  $v = U^T z$ .
- Return  $\underline{\Omega}_\mu^{-1} \underline{\Psi}^{-1} \left( \sum_{i=1}^N \underline{X}_i - \sum_{i=1}^N \underline{\Lambda} \underline{f}_i \right) + v$  (21).

### 3.5 Gibbs Sampler Pseudo-Code

- i) Choose scalar hyperparameters  $\Sigma_\mu, \Sigma_\lambda, \alpha$ , and  $\beta$ , and select  $q$ .
- ii) Initialise:
 
$$\begin{aligned} \underline{\mu}^{(0)} &\sim \text{MVN}_p(\underline{0}, \Sigma_\mu \mathcal{I}_p) \\ \underline{F}^{(0)} &\sim \text{MVN}_q(N, \underline{0}, \mathcal{I}_q) \\ \underline{\Lambda}^{(0)} &\sim \text{MVN}_q(N, \underline{0}, \Sigma_\lambda \mathcal{I}_q) \\ \underline{\Psi}^{(0)-1} &\sim \text{Ga}_p(\alpha/2, \beta/2) \end{aligned}$$
- iii) For  $t = 1, \dots, T$ , using the routines specified in (Equations (13), (16), (18) and (21)):

- a)  $\underline{\Omega}_\mu^{(t)} = (\Sigma_\mu \mathcal{I}_p)^{-1} + N \underline{\Psi}^{-1(t-1)}$
- b)  $\underline{\mu}^{(t)} \sim \text{MVN}_p \left( \underline{\Omega}_\mu^{-1(t)} \underline{\Psi}^{-1(t-1)} \left( \sum_{i=1}^N \underline{X}_i - \sum_{i=1}^N \underline{\Lambda}^{(t-1)} \underline{f}_i^{(t-1)} \right), \underline{\Omega}_\mu^{(t)} \right)$
- c)  $\underline{\Omega}_F^{(t)} = \mathcal{I}_q + \underline{\Lambda}^{T(t-1)} \underline{\Psi}^{-1(t-1)} \underline{\Lambda}^{(t-1)}$
- d)  $\underline{f}_i^{(t)} \sim \text{MVN}_q \left( \underline{\Omega}_F^{-1(t)} \underline{\Lambda}^{T(t-1)} \underline{\Psi}^{-1(t-1)} (\underline{X}_i - \underline{\mu}^{(t)}), \underline{\Omega}_F^{-1(t)} \right)$
- e) For  $j = 1, \dots, p$ 
  - $\underline{\Omega}_{\lambda_j}^{(t)} = (\Sigma_\lambda \mathcal{I}_q)^{-1} + \underline{\Psi}_{jj}^{-1(t-1)} \underline{F}^{T(t-1)} \underline{F}^{(t)}$
  - $\underline{\Lambda}_j^{(t)} \sim \text{MVN}_q \left( \underline{\Omega}_{\lambda_j}^{-1(t)} \underline{F}^{T(t)} \underline{\Psi}_{jj}^{-1(t-1)} (\underline{X}^j - \underline{\mu}_j^{(t)}), \underline{\Omega}_{\lambda_j}^{-1(t)} \right)$
- f)  $\underline{\Psi}^{-1(t)} \sim \text{Ga}_p \left( \frac{N + \alpha}{2}, \frac{S_{jj}^{(t)} + \beta}{2} \right)$

- iv) Disregard the first B burn-in iterations and thin every K-th iteration.



### 3.6 Issues Around Identifiability

Most covariance matrices  $\underline{\Sigma}$  cannot be uniquely factored as  $\underline{\Lambda}\underline{\Lambda}^T + \underline{\Psi}$  where  $q \ll p$ . Let  $\underline{T}$  be any  $q \times q$  orthogonal matrix such that  $\underline{T}\underline{T}^T = \underline{I}_q$ . Then:

$$\begin{aligned}\underline{x} - \underline{\mu} &= \underline{\Lambda}\underline{f} + \underline{\varepsilon} \\ &= \underline{\Lambda}\underline{T}\underline{T}^T\underline{f} + \underline{\varepsilon} \\ &= \underline{\Lambda}^*\underline{f}^* + \underline{\varepsilon}\end{aligned}$$

where  $\underline{\Lambda}^* = \underline{\Lambda}\underline{T}$  and  $\underline{f}^* = \underline{T}^T\underline{f}$ . It follows that  $E(\underline{f}^*) = \underline{0}$  and  $\text{Cov}(\underline{f}^*) = \underline{I}_q$ . Thus it is impossible, given the data  $\underline{x}$ , to distinguish between  $\underline{\Lambda}$  and  $\underline{\Lambda}^*$  since they both generate the same covariance matrix  $\underline{\Sigma}$ :

$$\begin{aligned}\underline{\Sigma} &= \underline{\Lambda}\underline{\Lambda}^T + \underline{\Psi} \\ &= \underline{\Lambda}\underline{T}\underline{T}^T\underline{\Lambda}^T + \underline{\Psi} \\ &= \underline{\Lambda}^*\underline{\Lambda}^{*T} + \underline{\Psi}\end{aligned}$$

However, we can solve this identifiability problem, using Procrustean methods, by mapping each iteration's loadings matrix to a common 'template' loadings matrix — which we have taken to be the loadings matrix at the burn-in iteration. This Procrustean map is a rotation only, i.e. translation, scaling, dilation, etc. are not applied. We then also apply that same rotation matrix at each iteration to each iteration of the matrix of factor scores. This amounts to *post-multiplying* the loadings matrix at each iteration by the Procrustes rotation matrix that maps to the loadings template, and also, by using the identity  $(AB)^T = B^T A^T$ , letting each iteration's scores matrix equal the transpose of the *post-multiplication* of that iteration's score matrix by that same rotation matrix.

## 4 Introducing the Shrinkage Prior

### 4.1 Multiplicative Gamma Process Shrinkage Priors

We now propose the multiplicative gamma process shrinkage prior of Bhattacharya & Dunson (2011) on the factor loadings which allows the introduction of infinitely many factors, with the loadings increasingly shrunk towards zero as the column index increases. Their prior is placed on a parameter expanded factor loadings matrix without imposing any restriction on the loading elements, thereby making the induced prior on the covariance matrix invariant to the ordering of the data. The Gibbs sampler can still be used due to the joint conjugacy property of this prior, which allows block updating of the loadings matrix. Furthermore, these authors propose that an adaptive Gibbs sampler be used for automatically truncating the infinite loading matrix, through selection of the number of important factors, to one having finite columns. This facilitates posterior computation while providing an accurate approximation to the infinite factor model.

The exact specification of this shrinkage-type prior has the degree of shrinkage increasing across the column index as follows:

$$\begin{aligned}\lambda_{jk} \mid \phi_{jk}, \tau_k &\sim N(0, \phi_{jk}^{-1} \tau_k^{-1}) \\ \text{s.t. } \underline{\lambda}_j \mid \underline{\phi}_j, \underline{\tau} &\sim \text{MVN}_{q^*}(\underline{0}, \underline{D}_j)\end{aligned}\tag{22}$$

$$\begin{aligned}\text{where } \underline{D}_j^{-1} &= \text{diag}(\phi_{j1}\tau_1, \dots, \phi_{jq^*}\tau_{q^*}) \\ \phi_{jk} &\sim \text{Ga}(\nu/2, \nu/2)\end{aligned}\tag{23}$$

$$\begin{aligned}\tau_k &= \prod_{h=1}^k \delta_h \\ \delta_1 &\sim \text{Ga}(\alpha_1, 1), \quad \delta_h \sim \text{Ga}(\alpha_2, 1), \quad h \geq 2\end{aligned}\tag{24}$$

where  $\delta_h$  ( $h = 1, \dots, \infty$ ) are independent,  $\tau_k$  is a *global* shrinkage parameter for the  $k$ -th column and the  $\phi_{jk}$ s are *local* shrinkage parameters for the elements in the  $k$ -th column. The  $\tau_k$ s are stochastically increasing under the restriction  $\alpha_2 > 1$ , which favours more shrinkage as the column index increases.

## 4.2 Deriving new MGP Full Conditionals

We propose a Gibbs sampler for posterior computation, much like the one above, after truncating the loadings matrix to have  $q^* \ll p$  columns. An adaptive strategy for inference on the truncation level  $q^*$  is described in (Section 4.3). For now, let's focus on the new full conditionals for the loadings matrix, global shrinkages, and local shrinkages which need to be derived in order to implement this. Once again, these parameters are initialised according to their priors. The other full conditionals are exactly as before, with just a small adjustment to the factor scores to allow for the truncation to  $q^*$  columns s.t.  $P(\underline{f}_i | \text{---}) \sim \text{MVN}_{q^*} \left( [\underline{\mathcal{I}}_{q^*} + \underline{\Lambda}_{q^*}^T \underline{\Psi}^{-1} \underline{\Lambda}_{q^*}]^{-1} \underline{\Lambda}^T \underline{\Psi}^{-1} \underline{X}_i^*, [\underline{\mathcal{I}}_{q^*} + \underline{\Lambda}_{q^*}^T \underline{\Psi}^{-1} \underline{\Lambda}_{q^*}]^{-1} \right)$

### 4.2.1 Loadings Matrix - $\underline{\Lambda}$

Incorporating the new prior (Equation (22)), and following the same steps as (Section 3.2) above, it is trivial to show that the  $\underline{\Lambda}_j$ s now have independent conditionally conjugate posteriors given by:

$$P(\underline{\Lambda}_j | \text{---}) \sim \text{MVN}_{q^*} \left( [\underline{D}_j^{-1} + \underline{\Psi}_{jj}^{-1} \underline{F}^T \underline{F}]^{-1} \underline{F}^T \underline{\Psi}_{jj}^{-1} \underline{X}^{j*}, [\underline{D}_j^{-1} + \underline{\Psi}_{jj}^{-1} \underline{F}^T \underline{F}]^{-1} \right) \quad (25)$$

However, we can reintroduce  $\underline{\mu}$  and save on computational time, as before, if we:

- Calculate  $\underline{\Omega}_{\lambda_j} = \underline{D}_j^{-1} + \underline{\Psi}_{jj}^{-1} \underline{F}^T \underline{F}$
- Compute the Cholesky Factorization  $\underline{\Omega}_{\lambda_j} = U^T U$ .
- Sample  $z \sim N(0, 1)$   $q$  times.
- Back-solve  $Uv = z$ .
- Compute  $\underline{\Omega}_{\lambda_j}^{-1}$  from  $U$ .
- Return  $\underline{\Omega}_{\lambda_j}^{-1} \underline{F}^T \underline{\Psi}_{jj}^{-1} (\underline{X}^j - \underline{\mu}_j) + v$  (26).

### 4.2.2 Local Shrinkage

Local Shrinkage –  $\phi_{jk}$  Using the conditional prior in (Equation (22)) and the prior for  $\phi_{jk}$  in (Equation (23)) we can derive the full conditional for the local shrinkage parameter as follows:

$$\begin{aligned} P(\phi_{jk} | \text{---}) &\propto P(\lambda_{jk} | \phi_{jk}, \tau_k) P(\phi_{jk}) \\ &\propto \frac{\phi_{jk}^{1/2} \tau_k^{1/2}}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \lambda_{jk}^2 \phi_{jk} \tau_k \right\} \phi_{jk}^{\nu/2-1} \exp \left\{ -\frac{\nu}{2} \phi_{jk} \right\} \\ &\propto \phi_{jk}^{1/2} \phi_{jk}^{\nu/2-1} \exp \left\{ \left( -\frac{1}{2} \lambda_{jk}^2 \tau_k - \frac{\nu}{2} \right) \phi_{jk} \right\} \\ &\propto \phi_{jk}^{\nu/2-1/2} \exp \left\{ -\frac{1}{2} (\nu + \lambda_{jk}^2 \tau_k) \phi_{jk} \right\} \end{aligned}$$

Thus the posterior distribution of each  $\phi_{jk}$  is given by:

$$P(\phi_{jk} | \text{---}) \sim \text{Ga} \left( \frac{\nu + 1}{2}, \frac{\nu + \tau_k \lambda_{jk}^2}{2} \right) \quad (27)$$

### 4.2.3 Global Shrinkage

Global Shrinkage –  $\tau_k$  Using the conditional prior in (Equation (22)) and the prior for  $\tau_k$  in (Equation (24)) we can derive the full conditional for the global shrinkage parameter, in three stages – first by deriving and sampling from  $P(\delta_1 | \text{---})$  &  $P(\delta_k | \text{---})$  for  $k \geq 2$ , as follows below — and then obtaining the product  $\tau_k = \prod_{h=1}^k \delta_h$  thereafter:

$$\begin{aligned}
P(\delta_1 | \text{---}) &\propto \prod_{j=1}^p \prod_{k=1}^{q^*} N(\lambda_{jk} | 0, \phi_{jk}^{-1} \tau_k^{-1}) \times \text{Ga}(\delta_1 | \alpha_1, 1) \\
&\propto \prod_{j=1}^p N(\lambda_{j1} | 0, \phi_{j1}^{-1} \tau_1^{-1}) \times \dots \times \prod_{j=1}^p N(\lambda_{jq^*} | 0, \phi_{jq^*}^{-1} \tau_{q^*}^{-1}) \times \text{Ga}(\delta_1 | \alpha_1, 1) \\
&\propto (\phi_{j1} \tau_1)^{p/2} \exp\left(-\frac{1}{2} \sum_{j=1}^p \lambda_{j1}^2 \phi_{j1} \tau_1\right) \times \dots \times (\phi_{jq^*} \tau_{q^*})^{p/2} \exp\left(-\frac{1}{2} \sum_{j=1}^p \lambda_{jq^*}^2 \phi_{jq^*} \tau_{q^*}\right) \\
&\quad \times \delta_1^{\alpha_1-1} \exp(-\delta_1) \\
&\propto (\phi_{j1} \delta_1)^{p/2} \exp\left(-\frac{1}{2} \sum_{j=1}^p \lambda_{j1}^2 \phi_{j1} \delta_1\right) \times \dots \times (\phi_{jq^*} \delta_1 \delta_2 \dots \delta_{q^*})^{p/2} \exp\left(-\frac{1}{2} \sum_{j=1}^p \lambda_{jq^*}^2 \phi_{jq^*} \delta_1 \delta_2 \dots \delta_{q^*}\right) \\
&\quad \times \delta_1^{\alpha_1-1} \exp(-\delta_1) \\
&\propto \delta_1^{pq^*/2 + \alpha_1 - 1} \exp\left(-\frac{\delta_1}{2} \left(\sum_{j=1}^p \lambda_{j1}^2 \phi_{j1} + \dots + \lambda_{jq^*}^2 \phi_{jq^*} \delta_2 \dots \delta_{q^*} + 2\right)\right) \\
&\propto \delta_1^{pq^*/2 + \alpha_1 - 1} \exp\left(-\frac{\delta_1}{2} \left(\sum_{h=1}^{q^*} \tau_h^{(1)} \sum_{j=1}^p \lambda_{jh}^2 \phi_{jh} + 2\right)\right)
\end{aligned}$$

$$\text{where } \tau_h^{(k)} = \prod_{t=1}^h \frac{\delta_t}{\delta_k} \text{ for } k = 1, \dots, q^* \quad (28)$$

$$\therefore P(\delta_1 | \text{---}) \sim \text{Ga}\left(\alpha_1 + \frac{pq^*}{2}, 1 + \frac{1}{2} \sum_{h=1}^{q^*} \tau_h^{(1)} \sum_{j=1}^p \lambda_{jh}^2 \phi_{jh}\right) \quad (29)$$

Similarly:

$$\begin{aligned}
P(\delta_k | \text{---}) &\propto \prod_{j=1}^p \prod_{k=1}^{q^*} N(\lambda_{jk} | 0, \phi_{jk}^{-1} \tau_k^{-1}) \times \text{Ga}(\delta_k | \alpha_2, 1) \\
&\propto (\phi_{j1} \delta_1)^{p/2} \exp\left(-\frac{1}{2} \sum_{j=1}^p \lambda_{j1}^2 \phi_{j1} \delta_1\right) \times \dots \times (\phi_{jq^*} \delta_1 \delta_2 \dots \delta_{q^*})^{p/2} \exp\left(-\frac{1}{2} \sum_{j=1}^p \lambda_{jq^*}^2 \phi_{jq^*} \delta_1 \delta_2 \dots \delta_{q^*}\right) \\
&\quad \times \delta_k^{\alpha_2-1} \exp(-\delta_k) \\
&\propto \delta_k^{p/2(q^*-k+1) + \alpha_2 - 1} \exp\left(-\frac{\delta_k}{2} \left(\sum_{h=k}^{q^*} \tau_h^{(k)} \sum_{j=1}^p \lambda_{jh}^2 \phi_{jh} + 2\right)\right) \\
&\therefore P(\delta_k | \text{---}) \sim \text{Ga}\left(\alpha_2 + \frac{p}{2}(q^* - k + 1), 1 + \frac{1}{2} \sum_{h=k}^{q^*} \tau_h^{(k)} \sum_{j=1}^p \lambda_{jh}^2 \phi_{jh}\right) \quad (30)
\end{aligned}$$

### 4.3 Adaptive Step

In practical situations, we expect to have relatively few important factors compared with the dimension  $p$  of the outcomes. The most common approach for selecting the number of factors relies on fitting the factor model for different choices of  $q^*$ , and then using the BIC or another criteria for model selection. This approach can be difficult to implement for large  $p$ , small  $N$  problems, and the BIC itself is not well

justified for factor models even for small to moderate  $p$ . However, the infinite factor model obviates the need for pre-specifying the number of factors, while the sparsity favouring prior on the loadings ensures that the effective number of factors would be small when the truth is sparse. However, we need a computational strategy for choosing an appropriate level of truncation  $q^*$ . We would like to strike a balance between missing important factors by choosing  $q^*$  too small and wasting computation on an overly conservative truncation level. One can think of  $q^*$  as the effective number of factors, so that the contribution from adding additional factors is negligible.

Starting with a conservative guess  $\tilde{q}$  of  $q^*$ , the posterior samples of  $\underline{\Lambda}_{\tilde{q}}$  from the Gibbs sampler contain information about the effective number of factors. At the  $t$ -th iteration of the Gibbs sampler, let  $m^{(t)}$  denote the number of columns in  $\underline{\Lambda}_{\tilde{q}}$  having all elements in a pre-specified small neighbourhood of zero. Intuitively,  $m^{(t)}$  of the factors have a negligible contribution at the  $t$ -th iteration. We then define  $q^{*(t)} = \tilde{q} - m^{(t)}$  to be the effective number of factors at iteration  $t$ .

It is typically necessary to choose a very conservative upper-bound to be assured that  $\tilde{q} \geq q^*$ , though this leads to wasted computational effort. Ideally, we would like to discard the redundant factors and continue sampling with a reduced number of loadings columns. We thereby save on computation by discarding unimportant factors. For this reason, the sampler described in (Section 3.5) above will be modified to an adaptive Gibbs sampler, which tunes the number of factors as the sampler progresses. We begin with a default value for  $\tilde{q}$  of the lesser of  $P$  and  $5 \ln(P)$ .

We adapt only after the burn-in period has elapsed, in order to ensure we're sampling from the true posterior distribution before truncating the loadings matrix. We adapt with probability  $p(t) = \exp(b_0 + b_1 t)$  at the  $t$ -th iteration after burn-in, with  $b_0, b_1$  chosen so that adaptation occurs around every 10 iterations at the beginning of the chain but decreases in frequency exponentially fast. We chose  $b_0$  and  $b_1$  in the adaptation probability  $p(t)$  as  $-0.1$  and  $-5 \times 10^{-5}$  respectively. We generate a sequence  $u_t$  of uniform random numbers between 0 and 1. At the  $t$ -th iteration, if  $u_t \leq p(t)$ , we monitor the columns in the loadings matrix having 75% of elements less than  $10^{-1}$  in magnitude. If the number of such columns drops to zero, we add a column to the loadings. Otherwise, we discard the redundant columns. The other parameters are also modified accordingly. When we add a factor, we sample parameters from their prior distributions to fill in additional columns, and otherwise retain parameters corresponding to the non-redundant columns. Letting  $\tilde{q}^{(t)}$  denote the truncation level at iteration  $t$  and  $q^{*(t)} = \tilde{q}^{(t)} - m^{(t)}$  denote the effective number of factors, we use the posterior mode or median of  $\{q^{*(t)}\}$  after burn-in as an estimate of  $q^*$  with credible intervals quantifying uncertainty.

## 5 Extension to Clustering Heterogeneous Data

### 5.1 Introducing Mixture Models

Marginally, (Equation 5) provides a parsimonious covariance matrix, s.t.  $\underline{X}_i \sim \text{MVN}_p(\underline{\mu}, \underline{\Lambda}\underline{\Lambda}^T + \underline{\Psi})$ . This allows us to exploit model-based clustering capabilities in high dimensional data settings. We can employ a(n) (in)finite mixture of factor analysis models whereby each of the  $G$  clusters is modelled using a cluster specific latent Gaussian model with covariance specified according to the form above. Let's now introduce some basic notation at this stage:

$$N = \sum_{g=1}^G n_g \quad \text{where } n_g \text{ is the size of the } g\text{-th group.}$$

$$P(\underline{X} | \underline{\gamma}) = \sum_{g=1}^G \pi_g P_g(\underline{X} | \underline{\theta}_g) \quad \text{where } \underline{\gamma} = (\underline{\theta}_1, \dots, \underline{\theta}_G, \pi_1, \dots, \pi_G) \text{ and the p.d.f } P_g \text{ is parametrized by } \underline{\theta}_g$$

$\pi_g$  are known as the *cluster mixing proportions*, which have the following properties

$$\begin{aligned} \pi_g &\geq 0 \quad \forall g = 1, \dots, G \\ \sum_{g=1}^G \pi_g &= 1 \end{aligned}$$

Introduce an additional latent indicator  $G$ -vector of *cluster labels* -  $\underline{z}_i$  - s.t.

$$z_{ig} = \begin{cases} 1 & \text{if } i \in g \\ 0 & \text{otherwise} \end{cases}$$

Therefore, if  $G = 3$  and observation  $i$  belongs to cluster 2,  $\underline{z}_i = (0, 1, 0)$ . Hence,

$$\begin{aligned} (\underline{X}_i | z_{ig} = 1) &\sim \text{MVN}_p(\underline{\mu}_g, \underline{\Lambda}_g \underline{\Lambda}_g^T + \underline{\Psi}_g) \\ \therefore P(\underline{X}_i) &= \sum_{g=1}^G \pi_g \text{MVN}_p(\underline{\mu}_g, \underline{\Lambda}_g \underline{\Lambda}_g^T + \underline{\Psi}_g) \end{aligned} \quad (31)$$

### 5.1.1 Decomposable Prior for $\underline{\gamma}$

The posterior distribution of  $\underline{\gamma}$  is

$$\begin{aligned} P(\underline{\gamma} | \underline{X}) &\propto P(\underline{\gamma}) \prod_{i=1}^N P(\underline{X}_i | \underline{\gamma}) \\ &\propto P(\underline{\gamma}) \prod_{i=1}^N \left( \sum_{g=1}^G \pi_g P_g(\underline{X}_i | \underline{\theta}_g) \right) \\ \therefore P(\underline{\gamma} | \underline{X}, \underline{z}) &\propto P(\underline{\gamma}) \prod_{g=1}^G \prod_{i: z_i=g} P_g(\underline{X}_i | \underline{\theta}_g) \end{aligned}$$

If, additionally,  $P(\underline{\gamma})$  can be decomposed into

$$\begin{aligned} P(\underline{\gamma}) &= P(\underline{\pi}) \prod_{g=1}^G P(\underline{\theta}_g), \\ P(\underline{\gamma} | \underline{X}, \underline{z}) &\propto P(\underline{\pi}) \prod_{g=1}^G \prod_{i: z_i=g} P(\underline{\theta}_g) P_g(\underline{X}_i | \underline{\theta}_g) \end{aligned} \quad (32)$$

## 5.2 Deriving Posterior Distributions

Attention now turns towards deriving full conditional distributions for the new parameter  $\underline{\pi}$ , as well as the latent variable  $\underline{z}$ , so that we can sample them for clustering purposes, by incorporating them into the Adaptive Gibbs Sampler framework described variously above.

- Component Parameters:

$$\begin{aligned} P(\underline{\theta}_g | \underline{\theta}_{-g}, \underline{X}, \underline{z}) &\equiv P(\underline{\theta}_g | \underline{X}, \underline{z}) \propto \prod_{i: z_i=g} P(\underline{\theta}_g) P_g(\underline{X}_i | \underline{\theta}_g) \\ \text{where } \underline{\theta}_{-1} &= (\underline{\theta}_1, \dots, \underline{\theta}_{g-1}, \underline{\theta}_{g+1}, \dots, \underline{\theta}_G) \end{aligned}$$

- Component Weights:

$$\begin{aligned} P(\underline{\pi} | \underline{X}, \underline{z}) &\equiv P(\underline{\pi} | \underline{z}) \propto P(\underline{\pi}) \prod_{g=1}^G \pi_g^{n_g} \\ \text{where } n_g &\text{ is the number of observations in group } g, \\ \text{since } P(\underline{z} | \underline{\pi}) &\sim \text{Mult}(N, \underline{\pi}) \end{aligned}$$

- Latent Classifier:

$$P(\underline{z} | \underline{X}, \underline{\gamma}) \propto \prod_{i=1}^N P(\underline{z}_i) P(\underline{X}_i | \underline{\theta}_i, \underline{z}_i)$$

### 5.2.1 Cluster Mixing Proportions – $\underline{\pi}$

Let the prior distribution of  $\underline{\pi}$  be Dirichlet with parameter  $\underline{\alpha}$  – a higher order generalisation of the Beta distribution.

$$\begin{aligned}
P(\underline{\pi}) &\propto \prod_{g=1}^G \underline{\pi}_g^{\alpha_g-1} \\
\therefore P(\underline{\pi} | \underline{z}, \underline{X}) &\propto \prod_{g=1}^G \underline{\pi}_g^{\alpha_g-1} \prod_{g=1}^G \underline{\pi}_g^{n_g} \\
&\propto \prod_{g=1}^G \underline{\pi}_g^{\alpha_g+n_g-1} \\
\text{s.t. } P(\underline{\pi} | \underline{z}, \underline{X}) &\sim \text{Dir}(\underline{\alpha} + \underline{n}) \\
&\text{where } \underline{n} = (n_1, \dots, n_G)
\end{aligned} \tag{33}$$

### 5.2.2 Cluster Labels – $\underline{z}_i$

$$\begin{aligned}
(\underline{z}_i | \underline{X}_i, \underline{\gamma}) &\sim \{1, \dots, G\}, \quad \forall i = 1, \dots, N \\
P(\underline{z}_i = g | \underline{X}_i, \underline{\gamma}) &= \frac{\underline{\pi}_g P(\underline{X}_i | \underline{\theta}_g)}{\sum_{g=1}^G \underline{\pi}_g P(\underline{X}_i | \underline{\theta}_g)}
\end{aligned}$$

### 5.2.3 MCMC Algorithm Pseudo-Code

1. Choose scalar hyperparameters as before.
2. Initialise:
  - (a) a
  - (b) b
  - (c) c
3. Start with some initialisation of the cluster labels  $\underline{z}^{(0)}$  from the  $\text{Mult}(N, \underline{\pi})$  prior.
4. Compute  $\underline{n}$ .
5. For  $g = 1, \dots, G$ , sample other parameters as before, but this time from their *group specific* full conditionals:
  - (a) a
  - (b) b
  - (c) c
6. Follow the adaptation procedure outlined in (Section 4.3).
7. Sample  $\underline{\pi}$  from  $\text{Dir}(\underline{\alpha} + \underline{n})$ .
8. For  $i = 1, \dots, N$ , sample  $\underline{z}_i$ .
9. Repeat steps 4–8 for  $t = 1, \dots, T$ .
10. Disregard the first B burn-in iterations and thin every K-th iteration.

### 5.3 Label Switching

When the main goal is identifying/interpreting the mixture components and/or clustering, then *label switching* should be treated. It is easy to see that  $P(\underline{X}|\underline{\gamma}) = P(\underline{X}|\underline{\tilde{\gamma}})$  where  $\underline{\tilde{\gamma}} = (\underline{\theta}_{j_1}, \dots, \underline{\theta}_{j_G}, \underline{\pi}_{j_1}, \dots, \underline{\pi}_{j_G})$  and  $j_1, \dots, j_G$  is any permutation of  $1, \dots, G$ . This type of finite mixture distribution nonidentifiability is caused by the invariance of a mixture distribution to relabelling the components: by interchanging the order of the components, the distributions induced by  $\underline{\gamma}$  and  $\underline{\tilde{\gamma}}$  are the same, although evidently the two parameters are distinct. For the finite mixture distribution as defined above with  $G$  components, there exist  $G!$  equivalent ways of arranging them.

## 6 References

- A. Bhattacharya & D. B. Dunson. *Sparse Bayesian infinite factor models*. *Biometrika*, 98(2): 291–306, 2011. ISSN 00063444. doi: 10.1093/biomet/asr013.
- S. Frühwirth-Schnatter. *Finite Mixture and Markov Switching Models*. Springer series in statistics, 2010. ISBN 9780387775005. doi: 10.1007/978-0-387-98135-2.
- G. J. McLachlan & D. Peel. *Finite mixture models*. Wiley series in probability and statistics. J. Wiley & Sons, New York, 2000. ISBN 0471006262. URL <http://opac.inria.fr/record=b1097397>.
- H. Rue & L. Held. *Gaussian Markov Random Fields: Theory and Applications*, volume 104 of *Monographs on Statistics and Applied Probability*. Chapman & Hall, London, 2005.