

# Infinite Mixtures of Infinite Factor Analysers

## Notes & Derivations

Keefe Murphy<sup>1, 2</sup>, Dr. Claire Gormley<sup>1, 2</sup>, Prof. Brendan Murphy<sup>1, 2</sup>, and Dr.  
Cinzia Viroli<sup>3</sup>

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

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

<sup>3</sup>Department of Statistical Sciences, University of Bologna

# 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 . . . . .	8
3.4	Reintroducing $\mu$ . . . . .	9
3.5	Gibbs Sampler Pseudo-Code . . . . .	9
3.6	Issues Around Identifiability . . . . .	10
<b>4</b>	<b>Introducing the Shrinkage Prior</b>	<b>11</b>
4.1	Multiplicative Gamma Process Shrinkage Priors . . . . .	11
4.2	Defining new MGP Full Conditionals . . . . .	11
4.2.1	Loadings Matrix . . . . .	11
4.2.2	Local Shrinkage . . . . .	12
4.2.3	Global Shrinkage . . . . .	12
4.3	Adaptive Step . . . . .	13
<b>5</b>	<b>Extension to Clustering Heterogeneous Data</b>	<b>14</b>
5.1	Introducing Mixture Models . . . . .	14
5.1.1	Decomposable Prior for $\gamma$ . . . . .	14
5.2	Deriving Posterior Distributions . . . . .	15
5.2.1	Cluster Mixing Proportions . . . . .	15
5.2.2	Latent Variables . . . . .	15
5.2.3	Mixtures of Infinite Factor Analyzers Pseudo-Code . . . . .	16
5.3	Label Switching . . . . .	17
5.4	Overfitting Mixtures . . . . .	17
5.5	Dirichlet Process Mixtures . . . . .	18
<b>6</b>	<b>References</b>	<b>18</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  $\Sigma$ . The factor model states that  $\underline{x}$  is linearly dependent upon a few ( $q \ll p$ ) unobservable random variables  $\underline{\eta}_1, \underline{\eta}_2, \dots, \underline{\eta}_q$  called *common factors* and  $p$  additional sources of variation  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_p$  called *specific factors*, for  $i = 1, \dots, N$  observations, s.t.

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

where	$\underline{x}_i \rightarrow (p \times 1)$	observation vector
	$\underline{\mu} \rightarrow (p \times 1)$	overall mean vector
	$\Lambda \rightarrow (p \times q)$	loadings matrix
	$\underline{\eta}_i \rightarrow (q \times 1)$	vector of factor scores for obs $i$
	$\underline{\varepsilon}_i \rightarrow (p \times 1)$	vector of errors for obs $i$

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

$$\left( \underline{x}_i - \underline{\mu} \right)_{(p \times 1)} = \underline{x}_{i(p \times 1)}^* = \Lambda_{(p \times q)} \underline{\eta}_{i(q \times 1)} + \underline{\varepsilon}_{i(p \times 1)} \quad (1.1)$$

## 1.2 Assumptions

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

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

$$\therefore \underline{\varepsilon}_i \sim \text{MVN}_p(0, \Psi) \quad (1.2)$$

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

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

$$\therefore \underline{\eta}_i \sim \text{MVN}_q(0, \mathcal{I}_q) \quad (1.3)$$

## 2 Bayesian Framework

### 2.1 Likelihood

$$\begin{aligned}
E(\underline{x}_i^*) &= E(\Lambda \underline{\eta}_i + \underline{\varepsilon}_i) \\
&= \Lambda E(\underline{\eta}_i) + E(\underline{\varepsilon}_i) \\
&= \underline{0} \\
\therefore \underline{x}_i^* &\sim \text{MVN}_p(\underline{0}, \Sigma)
\end{aligned} \tag{2.1}$$

$$\begin{aligned}
\text{Since } \underline{\varepsilon}_i &= \underline{x}_i^* - \Lambda \underline{\eta}_i, \\
\Sigma &= \text{Cov}(\underline{x}_i) \\
&= E\left[\left(\underline{x}_i - \underline{\mu}_i\right)\left(\underline{x}_i - \underline{\mu}_i\right)^T\right] \\
&= E\left[\underline{x}_i^* \underline{x}_i^{*T}\right] \\
&= E\left[\left(\Lambda \underline{\eta}_i + \underline{\varepsilon}_i\right)\left(\Lambda \underline{\eta}_i + \underline{\varepsilon}_i\right)^T\right] \\
&= E\left[\left(\Lambda \underline{\eta}_i\right) + \underline{\varepsilon}_i\left(\Lambda \underline{\eta}_i\right)^T + \left(\Lambda \underline{\eta}_i\right) \underline{\varepsilon}_i^T + \underline{\varepsilon}_i \underline{\varepsilon}_i^T\right] \\
&= \Lambda E\left(\underline{\eta}_i \underline{\eta}_i^T\right) \Lambda^T + E\left(\underline{\varepsilon}_i \underline{\eta}_i^T\right) \Lambda^T + \Lambda E\left(\underline{\eta}_i \underline{\varepsilon}_i^T\right) + E\left(\underline{\varepsilon}_i \underline{\varepsilon}_i^T\right) \\
&= \Lambda \Lambda^T + \Psi \\
\therefore \underline{x}_i^* &\sim \text{MVN}_p(\underline{0}, \Lambda \Lambda^T + \Psi)
\end{aligned} \tag{2.2}$$

$$\begin{aligned}
E(\underline{x}_i^* | \underline{\eta}_i) &= E(\Lambda \underline{\eta}_i + \underline{\varepsilon}_i | \underline{\eta}_i) \\
&= \Lambda E(\underline{\eta}_i | \underline{\eta}_i) + E(\underline{\varepsilon}_i | \underline{\eta}_i) \\
&= \Lambda \underline{\eta}_i \\
\text{Cov}(\underline{x}_i^* | \underline{\eta}_i) &= E\left[\left(\underline{x}_i^* - \Lambda \underline{\eta}_i\right)\left(\underline{x}_i^* - \Lambda \underline{\eta}_i\right)^T | \underline{\eta}_i\right] \\
&= E\left(\underline{\varepsilon}_i \underline{\varepsilon}_i^T | \underline{\eta}_i\right) \\
&= \Psi \\
\therefore \underline{x}_i^* | \underline{\eta}_i &\sim \text{MVN}_p(\Lambda \underline{\eta}_i, \Psi)
\end{aligned} \tag{2.3}$$

The density of the data is then given by:

$$P(\underline{x}_i^* | \underline{\eta}_i, \Lambda, \Psi) = (2\pi)^{-\frac{p}{2}} |\Psi|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \sum_{i=1}^N (\underline{x}_i^* - \Lambda \underline{\eta}_i)^T \Psi^{-1} (\underline{x}_i^* - \Lambda \underline{\eta}_i)\right) \quad (2.4)$$

$$\propto |\Psi|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \text{tr}\left[\Psi^{-1} (X - \eta \Lambda)^T (X - \eta \Lambda)\right]\right)$$

where  $\Lambda_{(p \times q)} = \begin{pmatrix} \lambda_{11} & \lambda_{12} & \dots & \lambda_{1q} \\ \lambda_{21} & \lambda_{22} & \dots & \lambda_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{p1} & \lambda_{p2} & \dots & \lambda_{pq} \end{pmatrix}$

&  $\eta_{(n \times q)} = \begin{pmatrix} \eta_{11} & \eta_{12} & \dots & \eta_{1q} \\ \eta_{21} & \eta_{22} & \dots & \eta_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ \eta_{n1} & \eta_{n2} & \dots & \eta_{nq} \end{pmatrix}$  &  $\underline{\eta}_i$  is a column vector containing the entries of row  $i$  of  $\eta$

## 2.2 Posterior Set-Up

$$\begin{aligned} \text{Likelihood} &= \prod_{i=1}^N P(\underline{x}_i^* | \theta) \\ &= \prod_{i=1}^N P(\underline{x}_i^* | \underline{\eta}_i, \Lambda, \Psi) \end{aligned}$$

$$\text{where } P(\underline{x}_i^* | \underline{\eta}_i, \Lambda, \Psi) \sim \text{MVN}_p(\Lambda \underline{\eta}_i, \Psi) \quad (2.5)$$

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

$$\text{Posterior} = \text{Likelihood} \times \text{Prior}$$

$$\begin{aligned} \therefore P(\eta, \Lambda, \Psi | X^*) &\propto \mathcal{L}(X^* | \eta, \Lambda, \Psi) P(\eta) P(\Lambda) P(\Psi) \\ &\propto \left[ \prod_{i=1}^N P(\underline{x}_i^* | \underline{\eta}_i, \Lambda, \Psi) \right] \left[ \prod_{i=1}^N P(\underline{\eta}_i) \right] \left[ \prod_{j=1}^p P(\underline{\Lambda}_j) \right] \left[ \prod_{j=1}^p P(\psi_j) \right] \end{aligned} \quad (2.6)$$

Later on, especially as we move into the mixture case, it will be necessary to undo the centering, thereby removing the  $*$  on  $\underline{x}_i^*$ , and reintroduce  $\underline{\mu}$ . This will necessitate multiplying the quantity in (2.6) by  $P(\underline{\mu})$ . 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 - $\underline{\eta}_i$

$$\begin{aligned} \underline{\eta}_i &\sim \text{MVN}_q(\underline{0}, \mathcal{I}_q) \\ &= (2\pi)^{-\frac{q}{2}} \exp\left(-\frac{1}{2} \underline{\eta}_i^T \underline{\eta}_i\right) \end{aligned} \quad (3.1)$$

To obtain the full conditional for  $\underline{\eta}_i$  we can multiply the likelihood by the prior in (3.1) s.t.

$$\begin{aligned}
P(\underline{\eta}_i | \underline{x}_i^*, \Lambda, \Psi) &\sim P(\underline{x}_i^* | \underline{\eta}_i, \Lambda, \Psi) P(\underline{\eta}_i) \\
&\propto \exp\left(-\frac{1}{2} \left[ (\underline{x}_i^* - \Lambda \underline{\eta}_i)^T \Psi^{-1} (\underline{x}_i^* - \Lambda \underline{\eta}_i) + \underline{\eta}_i^T \underline{\eta}_i \right]\right) \\
&\propto \exp\left(-\frac{1}{2} \left[ -\underline{x}_i^{*T} \Psi^{-1} \Lambda \underline{\eta}_i - (\Lambda \underline{\eta}_i)^T \Psi^{-1} \underline{x}_i^* + (\Lambda \underline{\eta}_i)^T \Psi^{-1} (\Lambda \underline{\eta}_i) + \underline{\eta}_i^T \underline{\eta}_i \right]\right) \\
&\propto \exp\left(-\frac{1}{2} \left\{ \underline{\eta}_i^T [\mathcal{I}_q + \Lambda^T \Psi^{-1} \Lambda] \underline{\eta}_i \right\} + \underline{x}_i^{*T} \Psi^{-1} \Lambda \underline{\eta}_i \right) \tag{3.2}
\end{aligned}$$

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

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

We can identify the  $\mu$  and  $\Sigma^{-1}$  terms from (3.2) above to yield

$$P(\underline{\eta}_i | \underline{x}_i^*, \Lambda, \Psi) \sim \text{MVN}_q\left([\mathcal{I}_q + \Lambda^T \Psi^{-1} \Lambda]^{-1} \Lambda^T \Psi^{-1} \underline{x}_i^*, [\mathcal{I}_q + \Lambda^T \Psi^{-1} \Lambda]^{-1}\right) \tag{3.3}$$

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  $\Omega_\eta = \mathcal{I}_q + \Lambda^T \Psi^{-1} \Lambda$
- Compute the Cholesky Factorization  $\Omega_\eta = U^T U$ .
- Sample  $\underline{z} \sim \text{MVN}_q(\underline{0}, \mathcal{I}_q)$   $N$  times.
- Backsolve  $U \underline{v} = \underline{z}^T$ .
- Compute  $\Omega_\eta^{-1}$  from  $U$ .
- Return  $\left(\Omega_\eta^{-1} \Lambda^T \Psi^{-1} (C_n \underline{\mu} X)^T + \underline{v}\right)^T$  (3.4)  
where  $C_n = \mathcal{I}_n - \frac{1}{n} \mathcal{O}$  and  $\mathcal{O}$  is an  $N \times N$  matrix of all 1's.

### 3.2 Loadings Matrix - $\Lambda$

A Gaussian distribution is a conjugate prior for  $\Lambda$ , implying an  $\text{MVN}_q$  distribution prior for each row  $\underline{\Lambda}_j$  of  $\Lambda$  s.t.  $\underline{\Lambda}_j \sim \text{MVN}_q(\underline{0}, \Sigma_\lambda)$  where  $\Sigma_\lambda$  is a diagonal covariance matrix. As above, we can

---

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

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

expect the result of the product of two  $\text{MVN}_q$  distributions to itself be distributed this way.

$$\begin{aligned}
P(\underline{\Lambda}_j | X^*, \eta, \Psi) &\sim P(X^* | \eta, \underline{\Lambda}_j, \Psi) P(\underline{\Lambda}_j | \Sigma_\lambda) \\
&\propto \exp \left( -\frac{1}{2} \sum_{i=1}^N \left( \underline{x}_i^* - \underline{\Lambda}_j \underline{\eta}_i \right)^T \psi_j^{-1} \left( \underline{x}_i^* - \underline{\Lambda}_j \underline{\eta}_i \right) \right) \exp \left( -\frac{1}{2} (\underline{\Lambda}_j^T \Sigma_\lambda^{-1} \underline{\Lambda}_j) \right) \\
&\propto \exp \left( -\frac{1}{2} \sum_{i=1}^N \left[ -2 \underline{x}_i^{*T} \psi_j^{-1} (\underline{\Lambda}_j \underline{\eta}_i) + (\underline{\Lambda}_j \underline{\eta}_i)^T \psi_j^{-1} (\underline{\Lambda}_j \underline{\eta}_i) + \underline{\Lambda}_j^T \Sigma_\lambda^{-1} \underline{\Lambda}_j \right] \right) \\
&\propto \exp \left( \underline{\Lambda}_j \psi_j^{-1} \sum_{i=1}^N \underline{x}_{ij}^{*T} \underline{\eta}_i - \frac{1}{2} \underline{\Lambda}_j^T \left[ \sum_{i=1}^N \psi_j^{-1} \underline{\eta}_i^T \underline{\eta}_i \right] \underline{\Lambda}_j - \frac{1}{2} \underline{\Lambda}_j^T \Sigma_\lambda^{-1} \underline{\Lambda}_j \right) \\
&\propto \exp \left( \underline{\Lambda}_j [\eta^T \psi_j^{-1} \underline{x}^{j*}] - \frac{1}{2} \underline{\Lambda}_j^T [\Sigma_\lambda^{-1} + \psi_j^{-1} \eta^T \eta] \underline{\Lambda}_j \right) \tag{3.5}
\end{aligned}$$

where  $\underline{x}^{j*}$  is an  $N$ -vector containing the elements of the  $j$ -th column of  $X^*$ .

$$\begin{aligned}
\therefore P(\underline{\Lambda}_j | X^*, \eta, \Psi) &\sim \text{MVN}_q \left( [\Sigma_\lambda^{-1} + \psi_j^{-1} \eta^T \eta]^{-1} \eta^T \psi_j^{-1} \underline{x}^{j*}, \right. \\
&\quad \left. [\Sigma_\lambda^{-1} + \psi_j^{-1} \eta^T \eta]^{-1} \right) \tag{3.6}
\end{aligned}$$

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

- Calculate  $\Omega_{\lambda_j} = \Sigma_{\lambda}^{-1} + \psi_j^{-1} \eta^T \eta$ .
  - Compute the Cholesky Factorization  $\Omega_{\lambda_j} = U^T U$ .
  - Sample  $\underline{z} \sim \text{MVN}_q(\underline{0}, \mathcal{I}_q)$ .
  - Back-solve  $U \underline{v} = \underline{z}$ .
  - Compute  $\Omega_{\lambda_j}^{-1}$  from  $U$ .
  - Return  $\Omega_{\lambda_j}^{-1} \eta^T \psi_j^{-1} (\underline{x}^j - \underline{1} \mu_j) + \underline{v}$   
where  $\underline{1}$  is an  $N$ -vector of all 1's.
- (3.7)

### 3.3 Uniquenesses - $\Psi$

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

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

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 get:

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

Since  $\Psi$  is a diagonal matrix:

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

This suggests the prior for  $\Psi^{-1}$  is a product of  $p$   $\text{Ga}(\alpha, \beta)$  distributions. We choose hyperparameters as per Frühwirth-Schnatter & Lopes (2010), by bounding each  $\psi_j$  away from zero in such a way that Heywood problems are avoided.

$$\begin{aligned} \therefore P(\Psi^{-1} | \alpha, \beta_j) &= \prod_{j=1}^p P(\psi_j^{-1} | \alpha, \beta_j) \\ &\propto \prod_{j=1}^p (\psi_j^{-1})^{\alpha-1} \exp(-\beta_j \psi_j^{-1}) \\ \therefore P(\Psi^{-1} | X^*, \eta, \Lambda) &\propto P(X^* | \eta, \Lambda) P(\Psi^{-1} | \alpha, \beta_j) \\ &\propto \prod_{j=1}^p (\psi_j^{-1})^{\frac{N}{2}} \exp\left(-\frac{\mathcal{S}_j^*}{2} \psi_j^{-1}\right) \prod_{j=1}^p (\psi_j^{-1})^{\alpha-1} \exp(-\beta_j \psi_j^{-1}) \\ &\propto \prod_{j=1}^p (\psi_j^{-1})^{\frac{N}{2} + \alpha - 1} \exp\left(-\left(\frac{\mathcal{S}_j^*}{2} + \beta_j\right) \psi_j^{-1}\right) \end{aligned} \quad (3.8)$$

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

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

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

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

$$P(\psi_j^{-1} | X, F, \Lambda) \sim \text{Ga}\left(\alpha + \frac{N}{2}, \beta_j + \frac{\mathcal{S}_j}{2}\right) \quad (3.9)$$



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

We've already seen from (3.4), (3.7) and (3.9) that reintroducing  $\mu$  to the other full conditionals is trivial. All that remains is to specify the conjugate Gaussian prior for  $\mu$  itself, and to derive its full conditional. This implies an  $\text{MVN}_p$  distribution prior s.t.  $\underline{\mu} \sim \text{MVN}_p(\tilde{\underline{\mu}}, \Sigma_\mu)$  where  $\Sigma_\mu$  is a diagonal covariance matrix, typically the diagonal of the data covariance matrix, and  $\tilde{\underline{\mu}}$  is a vector of prior mean means, typically the sample mean for each group. As above, we can expect the result of the product of two  $\text{MVN}_p$  distributions to itself be distributed this way.

$$\begin{aligned}
P(\underline{\mu} | X, \eta, \Psi, \Lambda) &\propto \exp\left(-\frac{1}{2} \sum_{i=1}^N \left(\underline{x}_i - \underline{\mu} - \Lambda \underline{\eta}_i\right)^T \Psi^{-1} \left(\underline{x}_i - \underline{\mu} - \Lambda \underline{\eta}_i\right)\right) \exp\left(-\frac{1}{2} (\underline{\mu} - \tilde{\underline{\mu}})^T \Sigma_\mu^{-1} (\underline{\mu} - \tilde{\underline{\mu}})\right) \\
&\propto \exp\left(-\frac{1}{2} \left( \sum_{i=1}^N \left[ -2 \underline{x}_i^T \Psi^{-1} \underline{\mu} + 2 \left(\Lambda \underline{\eta}_i\right)^T \Psi^{-1} \underline{\mu} + \underline{\mu}^T \Psi^{-1} \underline{\mu} \right] + \underline{\mu}^T \Sigma_\mu^{-1} \underline{\mu} - 2 \tilde{\underline{\mu}}^T \Sigma_\mu^{-1} \underline{\mu} \right)\right) \\
&\propto \exp\left( \sum_{i=1}^N \underline{x}_i^T \Psi^{-1} \underline{\mu} - \sum_{i=1}^N \left(\Lambda \underline{\eta}_i\right)^T \Psi^{-1} \underline{\mu} - \frac{1}{2} [\underline{\mu}^T (\Sigma_\mu^{-1} + N \Psi^{-1}) \underline{\mu}] + \tilde{\underline{\mu}}^T \Sigma_\mu^{-1} \underline{\mu} \right) \\
\therefore P(\underline{\mu} | X, \eta, \Psi, \Lambda) &\sim \text{MVN}_p \left( [\Sigma_\mu^{-1} + N \Psi^{-1}]^{-1} \left( \Psi^{-1} \left( \sum_{i=1}^N \underline{x}_i - \sum_{i=1}^N \Lambda \underline{\eta}_i \right) + \Sigma_\mu^{-1} \tilde{\underline{\mu}} \right), \right. \\
&\quad \left. [\Sigma_\mu^{-1} + N \Psi^{-1}]^{-1} \right) \tag{3.10}
\end{aligned}$$

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

- Calculate  $\Omega_\mu = \Sigma_\mu^{-1} + N \Psi^{-1}$ , which is a diagonal  $p \times p$  matrix.
- Invert  $\Omega_\mu$  by inverting its diagonal elements.
- $\Omega_\mu^{-1} = U^T U$  can be obtained by taking the square root of  $\Omega_\mu$  since this matrix is diagonal.
- Sample  $\underline{z} \sim \text{MVN}_p(\underline{0}, \mathcal{I}_p)$ .
- Compute  $\underline{v} = U^T \underline{z}$ .
- Return  $\Omega_\mu^{-1} \left( \Psi^{-1} \left( \sum_{i=1}^N \underline{x}_i - \sum_{i=1}^N \Lambda \underline{\eta}_i \right) + \Sigma_\mu^{-1} \tilde{\underline{\mu}} \right) + \underline{v}$  (3.11).

### 3.5 Gibbs Sampler Pseudo-Code

i) Choose hyperparameters  $\Sigma_\mu, \Sigma_\lambda, \alpha$ , and  $\beta$ , select  $q$  and initialise  $\tilde{\underline{\mu}}$ .

ii) Initialise:

$$\begin{aligned}
\underline{\mu}^{(0)} &\sim \text{MVN}_p(\tilde{\underline{\mu}}, \Sigma_\mu) \\
\underline{\eta}_i^{(0)} &\sim \text{MVN}_q(\underline{0}, \mathcal{I}_q) \quad \forall i = 1, \dots, N \\
\underline{\Lambda}_j^{(0)} &\sim \text{MVN}_q(\underline{0}, \Sigma_\lambda) \quad \forall j = 1, \dots, p \\
\psi_j^{-1(0)} &\sim \text{Ga}(\alpha, \beta_j) \quad \forall j = 1, \dots, p
\end{aligned}$$

iii) For  $t = 1, \dots, T$ , using the routines specified in (3.4), (3.7), (3.9) and (3.11):

$$\begin{aligned}
\text{a) } \Omega_{\mu}^{(t)} &= \Sigma_{\mu}^{-1} + N\Psi^{-1(t-1)} \\
\underline{\mu}^{(t)} &\sim \text{MVN}_p \left( \Omega_{\mu}^{-1(t)} \left( \Psi^{-1} \left( \sum_{i=1}^N \underline{x}_i - \sum_{i=1}^N \Lambda \underline{\eta}_i \right) + \Sigma_{\mu}^{-1} \tilde{\underline{\mu}} \right), \Omega_{\mu}^{-1(t)} \right) \\
\text{b) } \Omega_{\eta}^{(t)} &= \mathcal{I}_q + \Lambda^{T(t-1)} \Psi^{-1(t-1)} \Lambda^{(t-1)} \\
\underline{\eta}_i^{(t)} &\sim \text{MVN}_q \left( \Omega_{\eta}^{-1(t)} \Lambda^{T(t-1)} \Psi^{-1(t-1)} (\underline{x}_i - \underline{\mu}^{(t)}), \Omega_{\eta}^{-1(t)} \right) \\
\text{c) For } j &= 1, \dots, p \\
\bullet \Omega_{\lambda_j}^{(t)} &= \Sigma_{\lambda}^{-1} + \psi_j^{-1(t-1)} \eta^{T(t)} \eta^{(t)} \\
\underline{\Lambda}_j^{(t)} &\sim \text{MVN}_q \left( \Omega_{\lambda_j}^{-1(t)} \eta^{T(t)} \psi_j^{-1(t-1)} (\underline{x}^j - \underline{\mu}_j^{(t)}), \Omega_{\lambda_j}^{-1(t)} \right) \\
\bullet \psi_j^{-1(t)} &\sim \text{Ga} \left( \alpha + \frac{N}{2}, \beta_j + \frac{S_j^{(t)}}{2} \right)
\end{aligned}$$

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

v) Calculate the log-likelihood for each remaining sample. Then, using the largest value observed across these draws, BIC-MCMC, as defined by Frühwirth-Schnatter (2011), is determined by  $2 \ln \hat{\mathcal{L}} - k \ln(N)$ , where  $k = pq - \frac{q(q-1)}{2} + 2p$  is the effective number of parameters in the model. When choosing between competing models, the one with the highest BIC-MCMC is preferred. Alternatively, AIC-MCMC, or the BICM and AICM of Raftery et al. (2007) can be used.

### 3.6 Issues Around Identifiability

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

$$\begin{aligned}
\underline{x}_i - \underline{\mu} &= \Lambda \underline{\eta}_i + \varepsilon_i \\
&= \Lambda T T^T \underline{\eta}_i + \varepsilon_i \\
&= \Lambda^* \underline{\eta}_i^* + \varepsilon_i
\end{aligned}$$

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

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

However, we can address 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 end of the burn-in period. This Procrustean map is a rotation only, i.e. translation, scaling, dilation, etc. are not permitted. We then also apply that same rotation matrix at each iteration to each sample of the matrix of factor scores. This amounts to *post-multiplying* the loadings and factor score matrices at each iteration by the Procrustes rotation matrix that maps to that iteration's loadings template.

## 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 a close approximation of the infinite factor model.

The exact specification of this shrinkage-type prior allows the degree of shrinkage to increase across the column index as follows:

$$\lambda_{jk} | \phi_{jk}, \tau_k \sim N(0, \phi_{jk}^{-1} \tau_k^{-1})$$

$$\text{s.t. } \underline{\lambda}_j | \underline{\phi}_j, \underline{\tau} \sim \text{MVN}_{q^*}(\underline{0}, \underline{D}_j) \quad (4.1)$$

$$\text{where } \underline{D}_j^{-1} = \text{diag}(\phi_{j1}\tau_1, \dots, \phi_{jq^*}\tau_{q^*})$$

$$\phi_{jk} \sim \text{Ga}(\nu, \nu) \quad (4.2)$$

$$\tau_k = \prod_{h=1}^k \delta_h$$

$$\delta_1 \sim \text{Ga}(\alpha_1, \beta_1), \quad \delta_h \sim \text{Ga}(\alpha_2, \beta_2), \quad h \geq 2 \quad (4.3)$$

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  $\alpha_2 > \beta_2 + 1$ , which favours more shrinkage as the column index increases. Typically  $\beta_1 = \beta_2 = 1$ . However, we find it useful to consider an alternative reparameterisation of the local shrinkage prior via  $\phi_{jk} \sim \text{Ga}(\nu + 1, \nu)$ , s.t. the induced inverse-gamma prior on each  $\phi_{jk}^{-1}$  is non-informative in the sense that it has expectation 1.

### 4.2 Defining 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 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, i.e.  $P(\underline{\eta}_i | \text{---}) \sim \text{MVN}_{q^*}([ \mathcal{I}_{q^*} + \Lambda_{q^*}^T \Psi^{-1} \Lambda_{q^*} ]^{-1} \Lambda^T \Psi^{-1} \underline{x}_i - \underline{\mu}, [ \mathcal{I}_{q^*} + \Lambda_{q^*}^T \Psi^{-1} \Lambda_{q^*} ]^{-1})$

#### 4.2.1 Loadings Matrix - $\Lambda$

Incorporating the new prior (4.1), and following the same steps as 3.2 above, it's trivial to show that the  $\Lambda_j$ s now have independent conditionally conjugate posteriors given by:

$$P(\Lambda_j | \text{---}) \sim \text{MVN}_{q^*}([D_j^{-1} + \psi_j^{-1} \eta^T \eta]^{-1} \eta^T \psi_j^{-1} \underline{x}^{j^*}, [D_j^{-1} + \psi_j^{-1} \eta^T \eta]^{-1}) \quad (4.4)$$

However, we can reintroduce  $\underline{\mu}$  and save on computational time, as before, if we follow the routine given in (3.7), with  $\Omega_{\lambda_j} = D_j^{-1} + \psi_j^{-1} \eta^T \eta$ .

### 4.2.2 Local Shrinkage – $\phi_{jk}$

Using the conditional prior in (4.1) and the reparameterised version of (4.2), the prior for  $\phi_{jk}$ , we can derive the full conditional for the local shrinkage parameter as follows:

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

Thus the full conditional for each  $\phi_{jk}$  is given by:

$$P(\phi_{jk} | \text{---}) \sim \text{Ga} \left( \nu + \frac{3}{2}, \nu + \frac{\tau_k \lambda_{jk}^2}{2} \right) \quad (4.5)$$

### 4.2.3 Global Shrinkage – $\tau_k$

Using the conditional prior in (4.1) and the prior for  $\tau_k$  in (4.3) 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:

$$\begin{aligned} P(\delta_1 | \text{---}) &\propto \prod_{j=1}^p \prod_{k=1}^{q^*} N(\lambda_{jk} | 0, \phi_{jk}^{-1} \tau_k^{-1}) \times \text{Ga}(\delta_1 | \alpha_1, \beta_1) \\ &\propto \prod_{j=1}^p N(\lambda_{j1} | 0, \phi_{j1}^{-1} \tau_1^{-1}) \times \dots \times \prod_{j=1}^p N(\lambda_{jq^*} | 0, \phi_{jq^*}^{-1} \tau_{q^*}^{-1}) \times \text{Ga}(\delta_1 | \alpha_1, \beta_1) \\ &\propto (\phi_{j1} \tau_1)^{p/2} \exp \left( -\frac{1}{2} \sum_{j=1}^p \lambda_{j1}^2 \phi_{j1} \tau_1 \right) \times \dots \times (\phi_{jq^*} \tau_{q^*})^{p/2} \exp \left( -\frac{1}{2} \sum_{j=1}^p \lambda_{jq^*}^2 \phi_{jq^*} \tau_{q^*} \right) \\ &\quad \times \delta_1^{\alpha_1-1} \exp(-\beta_1 \delta_1) \\ &\propto (\phi_{j1} \delta_1)^{p/2} \exp \left( -\frac{1}{2} \sum_{j=1}^p \lambda_{j1}^2 \phi_{j1} \delta_1 \right) \times \dots \times (\phi_{jq^*} \delta_1 \delta_2 \dots \delta_{q^*})^{p/2} \exp \left( -\frac{1}{2} \sum_{j=1}^p \lambda_{jq^*}^2 \phi_{jq^*} \delta_1 \delta_2 \dots \delta_{q^*} \right) \\ &\quad \times \delta_1^{\alpha_1-1} \exp(-\beta_1 \delta_1) \\ &\propto \delta_1^{pq^*/2 + \alpha_1 - 1} \exp \left( -\frac{\delta_1}{2} \left( \sum_{j=1}^p \lambda_{j1}^2 \phi_{j1} + \dots + \lambda_{jq^*}^2 \phi_{jq^*} \delta_2 \dots \delta_{q^*} + 2\beta_1 \right) \right) \\ &\propto \delta_1^{pq^*/2 + \alpha_1 - 1} \exp \left( -\frac{\delta_1}{2} \left( \sum_{h=1}^{q^*} \tau_h^{(1)} \sum_{j=1}^p \lambda_{jh}^2 \phi_{jh} + 2\beta_1 \right) \right) \end{aligned}$$

$$\text{where } \tau_h^{(k)} = \prod_{t=1}^h \frac{\delta_t}{\delta_k} \text{ for } k = 1, \dots, q^* \quad (4.6)$$

$$\therefore P(\delta_1 | \text{---}) \sim \text{Ga} \left( \alpha_1 + \frac{pq^*}{2}, \beta_1 + \frac{1}{2} \sum_{h=1}^{q^*} \tau_h^{(1)} \sum_{j=1}^p \lambda_{jh}^2 \phi_{jh} \right) \quad (4.7)$$

$$\begin{aligned}
P(\delta_k | -) &\propto \prod_{j=1}^p \prod_{k=1}^{q^*} N(\lambda_{jk} | 0, \phi_{jk}^{-1} \tau_k^{-1}) \times \text{Ga}(\delta_k | \alpha_2, \beta_2) \\
&\propto (\phi_{j1} \delta_1)^{p/2} \exp\left(-\frac{1}{2} \sum_{j=1}^p \lambda_{j1}^2 \phi_{j1} \delta_1\right) \times \dots \times (\phi_{jq^*} \delta_1 \delta_2 \dots \delta_{q^*})^{p/2} \exp\left(-\frac{1}{2} \sum_{j=1}^p \lambda_{jq^*}^2 \phi_{jq^*} \delta_1 \delta_2 \dots \delta_{q^*}\right) \\
&\quad \times \delta_k^{\alpha_2-1} \exp(-\beta_2 \delta_k) \\
&\propto \delta_k^{p/2(q^*-k+1)+\alpha_2-1} \exp\left(-\frac{\delta_k}{2} \left(\sum_{h=k}^{q^*} \tau_h^{(k)} \sum_{j=1}^p \lambda_{jh}^2 \phi_{jh} + 2\beta_2\right)\right) \\
\therefore P(\delta_k | -) &\sim \text{Ga}\left(\alpha_2 + \frac{p}{2}(q^* - k + 1), \beta_2 + \frac{1}{2} \sum_{h=k}^{q^*} \tau_h^{(k)} \sum_{j=1}^p \lambda_{jh}^2 \phi_{jh}\right) \tag{4.8}
\end{aligned}$$

### 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 finite factor model for different choices of  $q^*$ , and then using the BIC, BIC-MCMC, or another model selection criterion. This approach can be difficult to implement for large  $p$ , small  $N$  problems, and the BIC itself isn't 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 high 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  $\Lambda_{\tilde{q}}$  from the Gibbs sampler contain information about the effective number of factors. At the  $t$ -th iteration, let  $m^{(t)}$  denote the number of columns in  $\Lambda_{\tilde{q}}$  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's 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 3.5 above is 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  $\min(\lfloor 3 \ln(p) \rfloor, p, N - 1)$ . 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 as  $-0.1$  and  $-5 \times 10^{-5}$  respectively. We generate a sequence  $u_t$  of uniform random numbers between 0 and 1. If  $u_t \leq p(t)$  at the  $t$ -th iteration, 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, an additional loadings column is added by simulating from the prior distribution. Otherwise redundant columns are discarded and parameters corresponding to the non-redundant columns are retained. The other parameters are also modified accordingly. 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. Thus a histogram approximation to the posterior for  $q^*$  is introduced and may be used to address the question about the number of latent factors.

## 5 Extension to Clustering Heterogeneous Data

### 5.1 Introducing Mixture Models

Marginally, 2.2 provides a parsimonious covariance matrix, i.e.  $\underline{x}_i | \theta \sim \text{MVN}_p(\underline{\mu}, \Lambda \Lambda^T + \Psi)$ . 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:

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

$$P(X | \gamma) = \sum_{g=1}^G \pi_g P_g(X | \theta_g) \quad \text{where } \gamma = (\theta_1, \dots, \theta_G, \pi_1, \dots, \pi_G), \quad (5.1)$$

and the p.d.f  $P_g$  is parametrized by  $\theta_g$ .

The *cluster mixing proportions* -  $\pi_1, \dots, \pi_G$  - have the following properties

$$\pi_g \geq 0 \quad \forall g = 1, \dots, G$$

$$\sum_{g=1}^G \pi_g = 1$$

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

$$\underline{x}_i | z_{ig} = 1 \sim \text{MVN}_p(\underline{\mu}_g, \Lambda_g \Lambda_g^T + \Psi_g)$$

$$\therefore P(\underline{x}_i) = \sum_{g=1}^G \pi_g \text{MVN}_p(\underline{\mu}_g, \Lambda_g \Lambda_g^T + \Psi_g) \quad (5.2)$$

#### 5.1.1 Decomposable Prior for $\gamma$

The posterior distribution of  $\gamma$  is

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

$$\propto P(\gamma) \prod_{i=1}^N \left( \sum_{g=1}^G \pi_g P_g(\underline{x}_i | \theta_g) \right)$$

$$\therefore P(\gamma | X, Z) \propto P(\gamma) \prod_{g=1}^G \prod_{i: z_{ig}=1} P_g(\underline{x}_i | \theta_g)$$

If,  $P(\gamma)$  can be decomposed into

$$P(\gamma) = P(\pi) \prod_{g=1}^G P(\theta_g), \text{ then}$$

$$P(\gamma | X, Z) \propto P(\pi) \prod_{g=1}^G \prod_{i: z_{ig}=1} P(\theta_g) P_g(\underline{x}_i | \theta_g) \quad (5.3)$$

## 5.2 Deriving Posterior Distributions

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

- Component Parameters –  $\theta_g$ :

$$P(\theta_g | \theta_{-g}, X, Z) \equiv P(\theta_g | X, Z) \propto \prod_{i: z_{ig}=1} P(\theta_g) P_g(\underline{x}_i | \theta_g)$$

$$\text{where } \theta_{-g} = (\theta_1, \dots, \theta_{g-1}, \theta_{g+1}, \dots, \theta_G)$$

- Cluster Mixing Proportions –  $\underline{\pi}$ :

$$P(\underline{\pi} | X, Z) \equiv P(\underline{\pi} | Z) \propto P(\underline{\pi}) \prod_{g=1}^G \pi_g^{n_g}$$

where  $n_g$  is the number of observations in group  $g$ ,

since  $P(\underline{z}_i | \underline{\pi}) \sim \text{Mult}(1, \underline{\pi})$

- Latent Variables –  $\underline{z}_i$ :

$$P(\underline{z}_i | \underline{x}_i, \gamma) \propto P(\underline{z}_i) P(\underline{x}_i | \theta_{i: z_{ig}=1}, \underline{z}_i)$$

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

Let the prior distribution of  $\underline{\pi}$  be Dirichlet with parameter  $\underline{\alpha}$  – a multivariate generalisation of the Beta distribution. Typically a symmetric uniform prior is chosen, whereby  $\alpha_g = 1 \forall g = 1, \dots, G$ .

$$\begin{aligned} P(\underline{\pi}) &\propto \prod_{g=1}^G \pi_g^{\alpha_g-1} \\ \therefore P(\underline{\pi} | Z, X) &\propto \prod_{g=1}^G \pi_g^{\alpha_g-1} \prod_{g=1}^G \pi_g^{n_g} \\ &\propto \prod_{g=1}^G \pi_g^{\alpha_g+n_g-1} \\ \text{i.e. } P(\underline{\pi} | Z, X) &\sim \text{Dir}(\underline{\alpha} + \underline{n}) \end{aligned} \tag{5.4}$$

where  $\underline{n} = (n_1, \dots, n_G)$

### 5.2.2 Latent Variables – $\underline{z}_i$

$\underline{z}_i | \underline{x}_i, \gamma \sim \text{Mult}(1, \underline{p})$ , where

$\underline{p} = (p_1, \dots, p_G)$ , and

$$\begin{aligned} p_g &= P(\underline{z}_{ig} = 1 | \underline{x}_i, \gamma) = \frac{\pi_g P(\underline{x}_i | \theta_g)}{\sum_{g=1}^G \pi_g P(\underline{x}_i | \theta_g)} = \frac{\pi_g f(\underline{x}_i | \underline{\mu}_g, \Lambda_g \Lambda_g^T + \Psi_g)}{\sum_{g=1}^G \pi_g f(\underline{x}_i | \underline{\mu}_g, \Lambda_g \Lambda_g^T + \Psi_g)} \\ &= \exp \left[ \log(\pi_g) + \log(f(\underline{x}_i | \underline{\mu}_g, \Lambda_g \Lambda_g^T + \Psi_g)) - \sum_{g=1}^G \left( \log(\pi_g) + \log(f(\underline{x}_i | \underline{\mu}_g, \Lambda_g \Lambda_g^T + \Psi_g)) \right) \right] \end{aligned} \tag{5.5}$$

### 5.2.3 Mixtures of Infinite Factor Analyzers Pseudo-Code

1. Choose scalar hyperparameters as before.
2. Start by initialising the cluster labels  $Z^{(0)}$ : simulate from the  $\text{Mult}(1, \underline{\pi})$  prior  $N$  times, or employ another clustering algorithm, such as K-Means. Compute  $\underline{n}$ , and  $\tilde{\underline{\mu}}_g$  for each group.
3. Initialise,  $\forall g = 1, \dots, G$ :

$$\begin{aligned}
\underline{\mu}_g^{(0)} &\sim \text{MVN}_p \left( \tilde{\underline{\mu}}_g, \Sigma_\mu \right) \\
\underline{\eta}_i^{(0)} &\sim \text{MVN}_{q_g^*} \left( \underline{0}, \mathcal{I}_{q_g^*} \right) \quad \forall i = 1, \dots, N \\
\underline{\Lambda}_{jg}^{(0)} &\sim \text{MVN}_{q_g^*} \left( \underline{0}, \Sigma_\lambda \right) \quad \forall j = 1, \dots, p \\
\psi_{jg}^{-1(0)} &\sim \text{Ga} \left( \alpha, \beta_j \right) \quad \forall j = 1, \dots, p \\
\phi_{jkg}^{(0)} &\sim \text{Ga} \left( \nu + 1, \nu \right) \quad \forall j = 1, \dots, p \quad \text{and} \quad k = 1, \dots, q_g^* \\
\delta_{1g}^{(0)} &\sim \text{Ga} \left( \alpha_1, \beta_1 \right), \quad \delta_{hg}^{(0)} \sim \text{Ga} \left( \alpha_2, \beta_2 \right), \quad h \geq 2 \\
\tau_{kg}^{(0)} &= \prod_{h=1}^k \delta_{hg}^{(0)} \quad \forall k = 1, \dots, q_g^*
\end{aligned}$$

4. For  $g = 1, \dots, G$ , sample other parameters as before, but this time from their *group specific* full conditionals:

$$\begin{aligned}
\text{a) } \Omega_{\mu_g}^{(t)} &= \Sigma_\mu^{-1} + n_g \Psi_g^{-1(t-1)} \\
\underline{\mu}_g^{(t)} &\sim \text{MVN}_p \left( \Omega_{\mu_g}^{-1(t)} \left( \Psi_g^{-1(t-1)} \left( \sum_{i:z_{ig}=1} \underline{x}_i - \sum_{i:z_{ig}=1} \Lambda_g^{(t-1)} \underline{\eta}_i^{(t-1)} \right) + \Sigma_\mu^{-1} \tilde{\underline{\mu}}_g \right), \Omega_{\mu_g}^{-1(t)} \right) \\
\text{b) } \Omega_{\eta_g}^{(t)} &= \mathcal{I}_{q_g^*} + \Lambda_g^{T(t-1)} \Psi_g^{-1(t-1)} \Lambda_g^{(t-1)} \\
\underline{\eta}_{i:z_{ig}=1}^{(t)} &\sim \text{MVN}_q \left( \Omega_{\eta_g}^{-1(t)} \Lambda_g^{T(t-1)} \Psi_g^{-1(t-1)} \left( \underline{x}_{i:z_{ig}=1} - \underline{\mu}_g^{(t)} \right), \Omega_{\eta_g}^{-1(t)} \right) \\
\text{c) For } j &= 1, \dots, p \\
&\bullet \Omega_{\lambda_{jg}}^{(t)} = \mathbf{D}_j^{-1} + \psi_{jg}^{-1(t-1)} \eta_{i:z_{ig}=1}^{T(t)} \eta_{i:z_{ig}=1}^{(t)} \\
&\quad \underline{\Lambda}_{jg}^{(t)} \sim \text{MVN}_{q_g^*} \left( \Omega_{\lambda_{jg}}^{-1(t)} \eta_{i:z_{ig}=1}^{T(t)} \psi_{jg}^{-1(t-1)} \left( \underline{x}_{i:z_{ig}=1}^j - \underline{\mu}_{jg}^{(t)} \right), \Omega_{\lambda_{jg}}^{-1(t)} \right) \\
&\bullet \psi_{jg}^{-1(t)} \sim \text{Ga} \left( \alpha + \frac{n_g}{2}, \beta_j + \frac{S_{jg}^{(t)}}{2} \right) \\
&\bullet \phi_{jkg}^{(t)} \sim \text{Ga} \left( \nu + \frac{3}{2}, \nu + \frac{\tau_{kg}^{(t-1)} \lambda_{jkg}^{2(t)}}{2} \right) \quad \forall k = 1, \dots, q_g^* \\
\text{d) } \delta_{1g}^{(t)} &\sim \text{Ga} \left( \alpha_1 + \frac{pq_g^*}{2}, \beta_1 + \frac{1}{2} \sum_{h=1}^{q_g^*} \tau_{hg}^{(1)(t-1)} \sum_{j=1}^p \lambda_{jhg}^{2(t)} \phi_{jhg}^{(t)} \right) \\
\delta_{hg}^{(t)} &\sim \text{Ga} \left( \alpha_2 + \frac{p}{2} (q_g^* - k + 1), \beta_2 + \frac{1}{2} \sum_{h=k}^{q_g^*} \tau_{hg}^{(k)(t)} \sum_{j=1}^p \lambda_{jhg}^{2(t-1)} \phi_{jhg}^{(t)} \right), \quad h \geq 2 \\
\tau_{kg}^{(t)} &= \prod_{h=1}^k \delta_{hg}^{(t)} \quad \forall k = 1, \dots, q_g^*
\end{aligned}$$

5. Re-compute  $\underline{n}$  and sample  $\underline{\pi}$  from  $\text{Dir}(\underline{\alpha} + \underline{n})$ .



6. For  $i = 1, \dots, N$ , sample  $z_i$  as outlined in (5.5).
7. Follow the adaptation procedure outlined in 4.3<sup>2</sup>.
8. Repeat steps 4–7 for  $t = 2, \dots, T$  using the current value for  $q_g^*$ .
9. Disregard the first B burn-in iterations and thin every K-th iteration <sup>3</sup>.

### 5.3 Label Switching

It's easy to see that  $P(X|\gamma) = P(X|\tilde{\gamma})$  where  $\tilde{\gamma} = (\theta_{j_1}, \dots, \theta_{j_G}, \pi_{j_1}, \dots, \pi_{j_G})$  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 mixture distributions to component relabelling: by interchanging the order of components, the distributions induced by  $\gamma$  and  $\tilde{\gamma}$  are the same, although evidently the two parameters are distinct. For finite mixture distribution as defined above with  $G$  components, there exist  $G!$  equivalent ways of arranging them. Generally as the Markov chain progresses, we will observe switches between these equivalent modes. When the main goal is identifying/interpreting mixture components &/or clustering, this *label switching* phenomenon needs to be addressed. The approach we adopt to do so is applied post-hoc, after the chain has finished running, and has the advantage of not involving loss functions based on sampled model parameters. We only require samples of  $Z$ , which are matched to a template vector of cluster labels at burnin using the cost-minimizing permutation suggested by the square assignment algorithm of Carpaneto & Toth (1980). This same permutation is applied to all other parameters which vary by group, namely the means, loadings, uniquenesses, and mixing proportions, prior to computing their posterior mean estimates.

### 5.4 Overfitting Mixtures

The need to choose the optimal number of latent factors in a mixture of factor analysers has been obviated using MIFA, but the issue of model choice is still not entirely resolved. Overfitting mixtures, along with Dirichlet Processes (5.5), are a means of extending the MIFA methodology in order to estimate  $G$  in a similarly choice-free manner. The prior in 5.2.1 plays an important role. This method approaches mixture model estimation by initially overfitting the number of clusters expected to be present, and specified conditions on the Dirichlet hyperparameter for the cluster mixing proportions encourage the emptying out of excess components in the posterior distribution.

To initialise the method, a conservatively high number of groups  $G^*$  is chosen, and fixed for the entire length of the MCMC chain. It's assumed that  $G^* > G$ . Each  $\alpha_g = 0.5/G^*$  is set small enough to favour empty groups a priori Ishwaran et al. (2001). The symmetric uniform prior  $\text{Dir}(1, \dots, 1)$  used previously is rather indifferent in this respect. The number of non-empty groups at each iteration  $G_0$  is recorded thusly:

$$G_0 = G^* - \sum_{g=1}^G \mathbb{1}(\sum_i z_{ig} = 0) \quad (5.6)$$

The true  $G$  is estimated by the  $G_0$  value visited most often by the sampler. Component specific inference is conducted only on the  $M_0$  samples corresponding to those visits.

---

<sup>2</sup>Our R-package also implements MFA, without the MGP shrinkage prior in 4.1 and adaptation.

<sup>3</sup>If using the MFA approach, one chooses between competing models according to the pair of G and Q values which optimise one of the model selection criteria outlined in 3.5. When using the MIFA approach, one chooses G using BICM or AICM only.

## 5.5 Dirichlet Process Mixtures

## 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.
- G. Carpaneto & P. Toth. *Algorithm 548 Solution of the assignment problem*. *ACM Transactions on Mathematical Software*, 6: 104–111, 1980.
- S. Frühwirth-Schnatter. *Finite Mixture and Markov Switching Models*. Springer series in statistics, 2010. ISBN 9780387775005. doi: 10.1007/978-0-387-98135-2.
- S. Frühwirth-Schnatter. *Dealing with label switching under model uncertainty*, pages 193–218. Mixture estimation and applications. Wiley, Chichester, 2011. ISBN ISBN-10: 11199938.
- S. Frühwirth-Schnatter & H. F. Lopes. Parsimonious bayesian factor analysis when the number of factors is unknown. (July 2015): 1–37, 2010.
- H. Ishwaran, L.F. James, & J. Sun. Bayesian model selection in finite mixtures by marginal density decompositions. *Journal of the American Statistical Association*, 96: 1316–1332, 2001. URL <http://EconPapers.repec.org/RePEc:bes:jnlasa:v:96:y:2001:m:december:p:1316-1332>.
- 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>.
- A. E. Raftery, M. Newton, P. N. Krivitsky, & J. M. Satagopan. Estimating the Integrated Likelihood via Posterior Simulation Using the Harmonic Mean Identity. *Bayesian Statistics*, (8): 1–45, 2007.
- H. Rue & L. Held. *Gaussian Markov Random Fields: Theory and Applications*, volume 104 of *Monographs on Statistics and Applied Probability*. Chapman & Hall, London, 2005.