

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

The Full Derivation of MSFA Model

Anonymous Author(s)
Affiliation
Address
email

1 Nonparametric Factor Analysis Model

Let $\mathbf{X} = [\mathbf{x}_{1:}, \dots, \mathbf{x}_{D:}]^T \in \mathbb{R}^{D \times N}$ be the data matrix in which \mathbf{x}_d denotes measured variable d and $\mathbf{x}_{:n}$ denotes data sample n . In standard factor analysis model, the data matrix \mathbf{X} can be decomposed into the product of the *factor loading matrix* $\mathbf{A} \in \mathbb{R}^{D \times K}$ and *factor matrix* $\mathbf{F} = [\mathbf{f}_{:1}, \dots, \mathbf{f}_{:N}] \in \mathbb{R}^{K \times N}$ plus noise matrix \mathbf{E} . The number of latent factors K is much smaller than D in general. The generative model of factor analysis is given by:

$$\begin{aligned} \mathbf{X} &= \mathbf{A}\mathbf{F} + \mathbf{E} \\ \mathbf{F} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \mathbf{E} &\sim \mathcal{N}(\mathbf{0}, \Psi) \end{aligned} \tag{1}$$

where the matrix $\mathbf{E} = [e_{:1}, \dots, e_{:N}]$ accounts for idiosyncratic noise. $\Psi = [\sigma_\epsilon^2 \mathbf{I}]$ is a diagonal matrix. We assume \mathbf{F} has a zero mean, unit variance Gaussian prior as used in standard factor analysis [?]. In nonparametric Bayesian factor analysis models [?, ?, ?, ?], the factor loading matrix \mathbf{A} can be factorized into the Hadamard product of a real-valued matrix \mathbf{G} and a binary-valued matrix \mathbf{Z} in further.

$$\mathbf{A} = \mathbf{G} \odot \mathbf{Z} \tag{2}$$

where \mathbf{G} and \mathbf{Z} are of the same size as \mathbf{A} . For each data sample n , the factor analysis model is given by: $\mathbf{x}_{:n} = (\mathbf{G} \odot \mathbf{Z})\mathbf{f}_{:n} + e_{:n}$. We assume \mathbf{G} has a Gaussian prior $\mathcal{N}(\mathbf{0}, \sigma_g^2 \mathbf{I})$ and \mathbf{Z} has a Dirichlet process prior $\mathcal{DP}(\beta, H_0)$ in which β is the concentration parameter. The graphical model for the nonparametric Bayesian factor analysis model and its equivalent representation are depicted in Figure (1)(2) respectively.

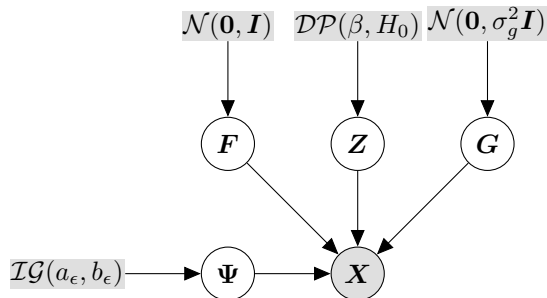


Figure 1: The Graphical Model for Nonparametric Bayesian Factor Analysis Model

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

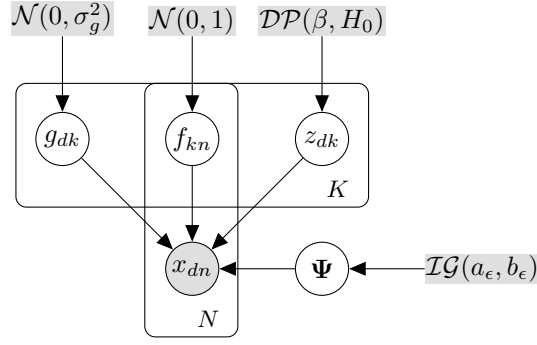


Figure 2: An Equivalent Graphical Representation for Nonparametric Bayesian Factor Analysis

We use Θ to denote model parameters $\{G, Z, F, \Psi\}$, and $x_{dn}, g_{dk}, z_{dk}, f_{kn}$ to denote the entry of $\mathbf{X}, \mathbf{G}, \mathbf{Z}, \mathbf{F}$ respectively.

The joint distribution of the data \mathbf{X} and model parameters Θ can be written as:

$$\begin{aligned}
p(\mathbf{X}, \Theta) &= \prod_{n=1}^N \prod_{d=1}^D \prod_{k=1}^K p(x_{dn}, g_{dk}, z_{dk}, f_{kn}, \sigma_\epsilon^2, \mu_g, \sigma_g^2, \beta, a_\epsilon, b_\epsilon) \\
&= \prod_{n=1}^N \prod_{d=1}^D \prod_{k=1}^K p(x_{dn} | g_{dk}, z_{dk}, f_{kn}, \sigma_\epsilon^2, \mu_g, \sigma_g^2, \beta, a_\epsilon, b_\epsilon) \\
&\quad \cdot \prod_{d=1}^D \prod_{k=1}^K p(g_{dk} | \mu_g, \sigma_g^2, z_{dk}, f_{kn}, \sigma_\epsilon^2, \beta, a_\epsilon, b_\epsilon) \\
&\quad \cdot \prod_{d=1}^D \prod_{k=1}^K p(z_{dk} | \mathbf{z}_{-dk}, f_{kn}, \sigma_\epsilon^2, \mu_g, \sigma_g^2, \beta, a_\epsilon, b_\epsilon) \\
&\quad \cdot \prod_{n=1}^N \prod_{k=1}^K p(f_{kn} | \sigma_\epsilon^2, \mu_g, \sigma_g^2, \beta, a_\epsilon, b_\epsilon) \\
&\quad \cdot p(\sigma_\epsilon^2 | a_\epsilon, b_\epsilon, \mu_g, \sigma_g^2, \beta) \\
&\quad \cdot p(\mu_g, \sigma_g^2, \beta, a_\epsilon, b_\epsilon) \\
&= \prod_{n=1}^N \prod_{d=1}^D \prod_{k=1}^K p(x_{dn} | g_{dk}, z_{dk}, f_{kn}, \sigma_\epsilon^2, \mathbf{z}_{-dk}, \mu_g, \sigma_g^2, \beta, a_\epsilon, b_\epsilon) \\
&\quad \cdot \prod_{d=1}^D \prod_{k=1}^K p(g_{dk} | \mu_g, \sigma_g^2, z_{dk}, \mathbf{z}_{-dk}, f_{kn}, \sigma_\epsilon^2, \beta, a_\epsilon, b_\epsilon) \\
&\quad \cdot \prod_{d=1}^D \prod_{k=1}^K p(z_{dk} | \mathbf{z}_{-dk}, \beta, f_{kn}, \sigma_\epsilon^2, \mu_g, \sigma_g^2, a_\epsilon, b_\epsilon) \\
&\quad \cdot \prod_{n=1}^N \prod_{k=1}^K p(f_{kn} | \sigma_\epsilon^2, \mu_g, \sigma_g^2, \beta, a_\epsilon, b_\epsilon) \\
&\quad \cdot p(\sigma_\epsilon^2 | a_\epsilon, b_\epsilon, \mu_g, \sigma_g^2, \beta) p(\mu_g, \sigma_g^2) p(\beta) p(a_\epsilon, b_\epsilon)
\end{aligned}$$

The variables with underline from line 8-12 can be removed from the conditional distribution according to the independencies among variables shown in Figure (1). We use a consistent concentration parameter β for all local factor analyzers as we assume the number of local latent factors is about the same.

The likelihood is given by:

$$\begin{aligned}
p(\mathbf{X}|\Theta) &= \prod_{n=1}^N \prod_{d=1}^D \prod_{k=1}^K p(x_{dn}|g_{dk}, z_{dk}, f_{kn}, \sigma_\epsilon^2) \\
&= \prod_{n=1}^N \prod_{d=1}^D \prod_{k=1}^K \left(\frac{1}{\sqrt{2\pi\sigma_\epsilon^2}} \exp \left[-\frac{(x_{dn} - g_{dk}f_{kn})^2}{2\sigma_\epsilon^2} \right] \right)^{z_{dk}} \\
&= \prod_{n=1}^N \prod_{d=1}^D \left(\frac{1}{\sqrt{2\pi\sigma_\epsilon^2}} \exp \left[-\frac{(x_{dn} - g_{d1}f_{1n})^2}{2\sigma_\epsilon^2} \right] \right)^{z_{d1}} \cdots \left(\frac{1}{\sqrt{2\pi\sigma_\epsilon^2}} \exp \left[-\frac{(x_{dn} - g_{dK}f_{Kn})^2}{2\sigma_\epsilon^2} \right] \right)^{z_{dK}} \\
&= \prod_{n=1}^N \prod_{d=1}^D \left(\frac{1}{\sqrt{2\pi\sigma_\epsilon^2}} \exp \left[-\frac{(x_{dn} - g_{d\hat{k}}f_{\hat{k}n})^2}{2\sigma_\epsilon^2} \right] \right) \\
&= \frac{1}{(2\pi\sigma_\epsilon^2)^{ND/2}} \exp \left[-\frac{\sum_{n=1, d=1}^{N, D} (x_{dn} - g_{d\hat{k}}f_{\hat{k}n})^2}{2\sigma_\epsilon^2} \right] \\
&= \frac{1}{(2\pi\sigma_\epsilon^2)^{ND/2}} \exp \left[-\frac{\text{tr}([\mathbf{X} - (\mathbf{G} \odot \mathbf{Z})\mathbf{F}])^T [\mathbf{X} - (\mathbf{G} \odot \mathbf{Z})\mathbf{F}]}{2\sigma_\epsilon^2} \right]
\end{aligned} \tag{3}$$

where we use \hat{k} to denote k when $z_{dk} = 1$.

We derive the posterior distribution of model parameters used in our Gibbs sampling in following subsections.

1.1 Sampling g_{dk}

We place a Gaussian prior $\mathcal{N}(\mu_g, \sigma_g^2)$ on each factor loading g_{dk} independently and let $\mu_g = 0$ as follows:

$$p(g_{dk}|\mu_g, \sigma_g^2) = \frac{1}{\sqrt{2\pi\sigma_g^2}} \exp \left[-\frac{(g_{dk} - \mu_g)^2}{2\sigma_g^2} \right] \tag{4}$$

Given the data and the other model parameters be fixed, we can derive the posterior distribution of g_{dk} as follows:

$$\begin{aligned}
p(g_{dk}|\mathbf{X}, -) &= \prod_{n=1}^N p(g_{dk}|x_{dn}, f_{kn}, z_{dk}, \mathbf{z}_{-dk}, \mu_g, \sigma_g^2, \sigma_\epsilon^2) \\
&= \prod_{n=1}^N \frac{p(x_{dn}, g_{dk}|z_{dk}, \mathbf{z}_{-dk}, f_{kn}, \sigma_\epsilon^2, \mu_g, \sigma_g^2)}{p(x_{dn}|z_{dk}, \mathbf{z}_{-dk}, f_{kn}, \sigma_\epsilon^2, \mu_g, \sigma_g^2)} \\
&= \prod_{n=1}^N \frac{p(x_{dn}|g_{dk}, z_{dk}, \mathbf{z}_{-dk}, f_{kn}, \sigma_\epsilon^2)p(g_{dk}|\mu_g, \sigma_g^2)}{p(x_{dn}|z_{dk}, \mathbf{z}_{-dk}, f_{kn}, \sigma_\epsilon^2, \mu_g, \sigma_g^2)}
\end{aligned}$$

162
163
164
165
166
167
168
169
170

$$\begin{aligned} &\propto \prod_{n=1}^N p(x_{dn}|g_{dk}, z_{dk}, \mathbf{z}_{-dk}, f_{kn}, \sigma_\epsilon^2) p(g_{dk}|\mu_g, \sigma_g^2) \\ & [\because p(x_{dn}|z_{dk}, \mathbf{z}_{-dk}, f_{kn}, \sigma_\epsilon^2, \mu_g, \sigma_g^2) \text{ is normalization constant given that the other parameters fixed.}] \\ &\propto \prod_{n=1}^N \left(\frac{1}{\sqrt{2\pi\sigma_\epsilon^2}} \exp \left[-\frac{(x_{dn} - g_{dk}f_{kn})^2}{2\sigma_\epsilon^2} \right] \right)^{z_{dk}} \frac{1}{\sqrt{2\pi\sigma_g^2}} \exp \left[-\frac{(g_{dk} - \mu_g)^2}{2\sigma_g^2} \right] \end{aligned}$$

171
172
173

where z_{dk} plays the role of feature (i.e.latent factor) selection as variable d can be associated to only one factor f_k : indicated by $z_{dk} = 1$ with the associating strength measured by g_{dk} . We use $-$ to denote the "rest" of the model parameters.

174
175

When $z_{dk} = 1$,

176
177
178
179
180
181
182
183
184
185

$$\begin{aligned} p(g_{dk}|X, -) &\propto \exp \left[-\frac{\sum_n (x_{dn} - g_{dk}f_{kn})^2}{2\sigma_\epsilon^2} \right] \exp \left[-\frac{(g_{dk} - \mu_g)^2}{2\sigma_g^2} \right] \\ & [\because \text{only the exponential terms involve } g_{dk}] \\ &\propto \exp \left[-\frac{\sum_n (x_{dn}^2 + g_{dk}^2 f_{kn}^2 - 2x_{dn}g_{dk}f_{kn})}{2\sigma_\epsilon^2} - \frac{g_{dk}^2 + \mu_g^2 - 2\mu_g g_{dk}}{2\sigma_g^2} \right] \\ &\propto \exp \left[-\frac{g_{dk}^2}{2} \left(\frac{\sum_n f_{kn}^2}{\sigma_\epsilon^2} + \frac{1}{\sigma_g^2} \right) + g_{dk} \left(\frac{\sum_n x_{dn}f_{kn}}{\sigma_\epsilon^2} + \frac{\mu_g}{\sigma_g^2} \right) - \left(\frac{\sum_n x_{dn}^2}{2\sigma_\epsilon^2} + \frac{\mu_g}{2\sigma_g^2} \right) \right] \end{aligned} \quad (5)$$

186

Since the product of two Gaussians is a Gaussian, we rewrite this in the form

187
188
189
190
191
192
193
194

$$\begin{aligned} p(g_{dk}|X, -) &\propto \exp \left[-\frac{g_{dk}^2}{2} \left(\frac{\sum_n f_{kn}^2}{\sigma_\epsilon^2} + \frac{1}{\sigma_g^2} \right) + g_{dk} \left(\frac{\sum_n x_{dn}f_{kn}}{\sigma_\epsilon^2} + \frac{\mu_g}{\sigma_g^2} \right) - \left(\frac{\sum_n x_{dn}^2}{2\sigma_\epsilon^2} + \frac{\mu_g}{2\sigma_g^2} \right) \right] \\ &= \exp \left[-\frac{\lambda_{dk}}{2} (g_{dk}^2 - 2g_{dk}\mu_{dk} + \mu_{dk}^2) \right] \\ &= \exp \left[-\frac{\lambda_{dk}}{2} (g_{dk} - \mu_{dk})^2 \right] \end{aligned} \quad (6)$$

195
196

Matching coefficient of g_{dk}^2 , we find λ_{dk} is given by

197
198
199

$$\lambda_{dk} = \frac{\sum_n f_{kn}^2}{\sigma_\epsilon^2} + \frac{1}{\sigma_g^2} \quad (7)$$

200
201

Matching coefficient of g_{dk} , we find μ_{dk} is given by

202
203
204

$$\mu_{dk} = \left(\frac{\sum_n f_{kn}^2}{\sigma_\epsilon^2} + \frac{1}{\sigma_g^2} \right)^{-1} \left(\frac{\sum_n x_{dn}f_{kn}}{\sigma_\epsilon^2} + \frac{\mu_g}{\sigma_g^2} \right) \quad (8)$$

205

206
207
208

Hence, given the data and the other parameters be fixed, we can sample $g_{dk} \sim \mathcal{N}(\mu_{dk}, \lambda_{dk})$ in which μ_{dk}, λ_{dk} , that are the mean and precision parameters of the posterior distribution of g_{dk} , can be calculated according to:

209
210
211
212
213
214
215

$$\begin{aligned} \lambda_{dk} &= \begin{cases} \frac{\sum_n f_{kn}^2}{\sigma_\epsilon^2} + \frac{1}{\sigma_g^2} & \text{if } z_{dk} = 1 \\ \frac{1}{\sigma_g^2} & \text{if } z_{dk} = 0 \end{cases} \\ \mu_{dk} &= \begin{cases} \left(\frac{\sum_n f_{kn}^2}{\sigma_\epsilon^2} + \frac{1}{\sigma_g^2} \right)^{-1} \left(\frac{\sum_n x_{dn}f_{kn}}{\sigma_\epsilon^2} + \frac{\mu_g}{\sigma_g^2} \right) & \text{if } z_{dk} = 1 \\ \mu_g & \text{if } z_{dk} = 0 \end{cases} \end{aligned} \quad (9)$$

216
217

1.2 Sampling z_{dk}

218
219
220
221

Our goal is to partition the measure variables into K clusters but K is unbounded a priori. The Dirichlet process (DP) defines a distribution over clustering where the number of clusters does not need to be specified a priori. The mixing proportions, $\omega_k = \frac{\sum_{d=1}^D z_{dk}}{D}$, are given a symmetric Dirichlet prior with the concentration parameters β :

222
223
224
225

$$p(\omega_1, \dots, \omega_K | \beta) \sim \text{Dirichlet}(\beta/K, \dots, \beta/K) = \frac{\Gamma(\beta)}{\Gamma(\beta/K)^K} \prod_{k=1}^K \omega_k^{\beta/K-1} \quad (10)$$

226
227
228
229
230
231

where the mixing proportions must be positive and sum to one. We use $z_{dk} = 1$ to indicate that the measured variable d is partitioned into cluster k . As we enforce the disjoint partition over variables, one variable can be grouped into only one cluster and thus $\sum_{k=1}^K z_{dk} = 1$. Given the mixing proportions, the prior for the occupation numbers, $m_k = \sum_{d=1}^D z_{dk}$, is multinomial and the joint distribution of the indicators becomes:

232
233
234
235

$$p(\{z_{dk}\}_{d=1, k=1}^{D, K} | \omega_1, \dots, \omega_K) = \prod_{k=1}^K \omega_k^{m_k} \quad (11)$$

236
237
238

Using the standard Dirichlet integral, we can integrate out the mixing proportions and write the prior directly in terms of the indicators as follows:

239
240
241

$$\begin{aligned} p(\{z_{dk}\}_{d=1, k=1}^{D, K} | \beta) &= \int p(\{z_{dk}\}_{d=1, k=1}^{D, K} | \omega_1, \dots, \omega_K) p(\omega_1, \dots, \omega_K | \beta) d\omega_1 \dots d\omega_K \\ &= \frac{\Gamma(\beta)}{\Gamma(\beta/K)^K} \int \prod_{k=1}^K \omega_k^{m_k + \beta/K - 1} d\omega_k \\ &= \frac{\Gamma(\beta)}{\Gamma(\beta/K)^K} \int \frac{\prod_k \Gamma(m_k + \beta/K)}{\Gamma(\sum_k m_k + \beta)} \frac{\Gamma(\sum_k m_k + \beta)}{\prod_k \Gamma(m_k + \beta/K)} \prod_{k=1}^K \omega_k^{m_k + \beta/K - 1} d\omega_k \\ &= \frac{\Gamma(\beta)}{\Gamma(\beta/K)^K} \frac{\prod_k \Gamma(m_k + \beta/K)}{\Gamma(\sum_k m_k + \beta)} \\ &= \frac{\Gamma(\beta)}{\Gamma(D + \beta)} \prod_{k=1}^K \frac{\Gamma(m_k + \beta/K)}{\Gamma(\beta/K)} \\ &[\because \int \frac{\Gamma(\sum_k m_k + \beta)}{\prod_k \Gamma(m_k + \beta/K)} \prod_{k=1}^K \omega_k^{m_k + \beta/K - 1} d\omega_k = 1] \end{aligned} \quad (12)$$

242
243
244
245
246
247
248
249
250
251
252
253

254
255
256

We use \mathbf{z}_{-dk} denotes all the indicators except $\mathbf{z}_d = [z_{d1}, \dots, z_{dK}]$. By keeping \mathbf{z}_{-dk} to be fixed, we obtain the conditional prior for the single indicator z_{dk} :

257
258
259
260
261
262

$$p(z_{dk} | \mathbf{z}_{-dk}, \beta) = \begin{cases} \frac{m_{-dk} + \beta/K}{D - 1 + \beta} & \text{if } z_{dk} = 1 \\ \frac{D - 1 + \beta - m_{-dk} - \beta/K}{D - 1 + \beta} & \text{if } z_{dk} = 0 \end{cases} \quad (13)$$

263
264
265

where the subscript $-d$ denotes all indexes except d and m_{-dk} is the number of variables, excluding \mathbf{x}_d , that are grouped into cluster k .

266

By taking the limit $K \rightarrow \infty$, we get the conditional prior for a single indicator as:

267
268
269

$$p(z_{dk} | \mathbf{z}_{-dk}, \beta) = \begin{cases} \frac{m_{-dk}}{D - 1 + \beta} & \text{if } z_{dk} = 1 \\ \frac{D - 1 + \beta - m_{-dk}}{D - 1 + \beta} & \text{if } z_{dk} = 0 \end{cases} \quad (14)$$

Given the data \mathbf{X} and model parameters be fixed excluding z_{dk} , we can calculate the posterior distribution of z_{dk} as:

$$\begin{aligned}
p(z_{dk} = 1 | \mathbf{X}, -) &= \prod_{n=1}^N p(z_{dk} = 1 | x_{dn}, g_{dk}, f_{kn}, \mathbf{z}_{-dk}, \beta) \\
&= \prod_{n=1}^N \frac{p(z_{dk} = 1, x_{dn} | g_{dk}, f_{kn}, \mathbf{z}_{-dk}, \beta)}{p(x_{dn} | g_{dk}, f_{kn}, \mathbf{z}_{-dk}, \beta)} \\
&\propto \prod_{n=1}^N p(z_{dk} = 1, x_{dn} | g_{dk}, f_{kn}, \mathbf{z}_{-dk}, \beta) \\
&= \prod_{n=1}^N p(x_{dn} | g_{dk}, z_{dk} = 1, f_{kn}, \mathbf{z}_{-dk}, \beta) p(z_{dk} = 1 | g_{dk}, f_{kn}, \mathbf{z}_{-dk}, \beta) \\
&= \prod_{n=1}^N p(x_{dn} | g_{dk}, z_{dk} = 1, f_{kn}, \mathbf{z}_{-dk}, \beta) p(z_{dk} = 1 | g_{dk}, f_{kn}, \mathbf{z}_{-dk}, \beta) \\
&= \prod_{n=1}^N p(x_{dn} | g_{dk}, z_{dk} = 1, f_{kn}, \beta) p(z_{dk} = 1 | \mathbf{z}_{-dk}, \beta) \\
&= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_\epsilon^2}} \exp\left[-\frac{(x_{dn} - g_{dk}f_{kn})^2}{2\sigma_\epsilon^2}\right] \frac{m_{-dk}}{D-1+\beta} \\
&= \frac{1}{(2\pi\sigma_\epsilon^2)^{N/2}} \exp\left[-\frac{\sum_{n=1}^N (x_{dn} - g_{dk}f_{kn})^2}{2\sigma_\epsilon^2}\right] \frac{m_{-dk}}{D-1+\beta} \\
&\propto \exp\left[-\frac{\text{tr}([\mathbf{x}_d - g_{dk}\mathbf{f}_k:]^T [\mathbf{x}_d - g_{dk}\mathbf{f}_k:])}{2\sigma_\epsilon^2}\right] \frac{m_{-dk}}{D-1+\beta}
\end{aligned}$$

Alternatively, we can integrate out g_{dk} to sample z_{dk} as:

$$\begin{aligned}
p(z_{dk} = 1 | \mathbf{X}, -) &= \int p(z_{dk} = 1, g_{dk} | \mathbf{X}, \Theta) dg_{dk} \\
&= \int \prod_{n=1}^N p(z_{dk} = 1, g_{dk} | x_{dn}, f_{kn}, \mathbf{z}_{-dk}, \beta, \mu_g, \sigma_g^2) dg_{dk} \\
&= \int \prod_{n=1}^N \frac{p(z_{dk} = 1, g_{dk}, x_{dn} | f_{kn}, \mathbf{z}_{-dk}, \beta, \mu_g, \sigma_g^2)}{p(x_{dn} | f_{kn}, \mathbf{z}_{-dk}, \beta, \mu_g, \sigma_g^2)} dg_{dk} \\
&\propto \int \prod_{n=1}^N p(z_{dk} = 1, g_{dk}, x_{dn} | f_{kn}, \mathbf{z}_{-dk}, \beta, \mu_g, \sigma_g^2) dg_{dk} \\
&= \int \prod_{n=1}^N p(x_{dn} | g_{dk}, z_{dk} = 1, f_{kn}, \mathbf{z}_{-dk}, \beta, \mu_g, \sigma_g^2) \\
&\quad \times p(z_{dk} = 1 | g_{dk}, f_{kn}, \mathbf{z}_{-dk}, \beta, \mu_g, \sigma_g^2) p(g_{dk} | \mu_g, \sigma_g^2) dg_{dk} \\
&= \int \prod_{n=1}^N p(x_{dn} | g_{dk}, z_{dk} = 1, f_{kn}) p(z_{dk} = 1 | \mathbf{z}_{-dk}, \beta) p(g_{dk} | \mu_g, \sigma_g^2) dg_{dk} \\
&= p(z_{dk} = 1 | \mathbf{z}_{-dk}, \beta) \int \prod_{n=1}^N p(x_{dn} | g_{dk}, z_{dk} = 1, f_{kn}) p(g_{dk} | \mu_g, \sigma_g^2) dg_{dk}
\end{aligned}$$

$$\begin{aligned}
&= p(z_{dk} = 1 | \mathbf{z}_{-dk}, \beta) \int \frac{1}{(2\pi\sigma_\epsilon^2)^{\frac{N}{2}}} \exp \left[-\frac{\sum_n (x_{dn} - g_{dk} f_{kn})^2}{2\sigma_\epsilon^2} \right] \frac{1}{(2\pi\sigma_g^2)^{\frac{1}{2}}} \exp \left[-\frac{(g_{dk} - \mu_g)^2}{2\sigma_g^2} \right] dg_{dk} \\
&= p(z_{dk} = 1 | \mathbf{z}_{-dk}, \beta) \int \frac{1}{(2\pi\sigma_\epsilon^2)^{\frac{N}{2}}} \exp \left[-\frac{\sum_n (x_{dn}^2 + g_{dk}^2 f_{kn}^2 - 2x_{dn}g_{dk}f_{kn})}{2\sigma_\epsilon^2} \right] \\
&\quad \times \frac{1}{(2\pi\sigma_g^2)^{\frac{1}{2}}} \exp \left[-\frac{g_{dk}^2 + \mu_g^2 - 2\mu_g g_{dk}}{2\sigma_g^2} \right] dg_{dk} \\
&= p(z_{dk} = 1 | \mathbf{z}_{-dk}, \beta) \int \frac{1}{(2\pi\sigma_\epsilon^2)^{\frac{N}{2}}} \exp \left[-\frac{g_{dk}^2}{2} \left(\frac{\sum_n f_{kn}^2}{\sigma_\epsilon^2} + \frac{1}{\sigma_g^2} \right) + g_{dk} \left(\frac{\sum_n x_{dn} f_{kn}}{\sigma_\epsilon^2} + \frac{\mu_g}{\sigma_g^2} \right) \right] \\
&\quad \times \frac{1}{(2\pi\sigma_g^2)^{\frac{1}{2}}} \exp \left[-\frac{\sum_n x_{dn}^2}{2\sigma_\epsilon^2} - \frac{\mu_g}{2\sigma_g^2} \right] dg_{dk} \\
&\quad \text{[Recall equation 6]} \\
&= p(z_{dk} = 1 | \mathbf{z}_{-dk}, \beta) \int \frac{1}{(2\pi\sigma_\epsilon^2)^{\frac{N}{2}}} \exp \left[-\frac{g_{dk}^2}{2} \lambda_{dk} + g_{dk} \lambda_{dk} \mu_{dk} + \frac{\lambda_{dk} \mu_{dk}^2}{2} - \frac{\lambda_{dk} \mu_{dk}^2}{2} \right] \\
&\quad \times \frac{1}{(2\pi\sigma_g^2)^{\frac{1}{2}}} \exp \left[-\frac{\sum_n x_{dn}^2}{2\sigma_\epsilon^2} - \frac{\mu_g}{2\sigma_g^2} \right] dg_{dk} \\
&= p(z_{dk} = 1 | \mathbf{z}_{-dk}, \beta) \int \frac{1}{(2\pi\sigma_\epsilon^2)^{\frac{N}{2}}} \exp \left[-\frac{\lambda_{dk}}{2} (g_{dk} - \mu_{dk})^2 \right] \exp \left[\frac{\lambda_{dk} \mu_{dk}^2}{2} \right] \\
&\quad \times \frac{1}{(2\pi\sigma_g^2)^{\frac{1}{2}}} \exp \left[-\frac{\sum_n x_{dn}^2}{2\sigma_\epsilon^2} - \frac{\mu_g}{2\sigma_g^2} \right] dg_{dk} \\
&= p(z_{dk} = 1 | \mathbf{z}_{-dk}, \beta) \int \exp \left[-\frac{\lambda_{dk}}{2} (g_{dk} - \mu_{dk})^2 \right] dg_{dk} \frac{1}{(2\pi\sigma_\epsilon^2)^{\frac{N}{2}}} \exp \left[\frac{\lambda_{dk} \mu_{dk}^2}{2} \right] \\
&\quad \times \frac{1}{(2\pi\sigma_g^2)^{\frac{1}{2}}} \exp \left[-\frac{\sum_n x_{dn}^2}{2\sigma_\epsilon^2} - \frac{\mu_g}{2\sigma_g^2} \right] \\
&= p(z_{dk} = 1 | \mathbf{z}_{-dk}, \beta) \frac{(2\pi)^{\frac{1}{2}}}{\lambda_{dk}^{\frac{1}{2}}} \frac{1}{(2\pi\sigma_\epsilon^2)^{\frac{N}{2}}} \exp \left[\frac{\lambda_{dk} \mu_{dk}^2}{2} \right] \frac{1}{(2\pi\sigma_g^2)^{\frac{1}{2}}} \exp \left[-\frac{\sum_n x_{dn}^2}{2\sigma_\epsilon^2} - \frac{\mu_g}{2\sigma_g^2} \right] \\
&= p(z_{dk} = 1 | \mathbf{z}_{-dk}, \beta) \frac{1}{\lambda_{dk}^{\frac{1}{2}}} \frac{1}{(2\pi\sigma_\epsilon^2)^{\frac{N}{2}}} \exp \left[\frac{\lambda_{dk} \mu_{dk}^2}{2} \right] \frac{1}{\sigma_g} \exp \left[-\frac{\sum_n x_{dn}^2}{2\sigma_\epsilon^2} - \frac{\mu_g}{\sigma_g^2} \right] \\
&\propto p(z_{dk} = 1 | \mathbf{z}_{-dk}, \beta) \frac{1}{\lambda_{dk}^{\frac{1}{2}}} \exp \left[\frac{\lambda_{dk} \mu_{dk}^2}{2} \right] \\
&\propto \frac{m_{-dk}}{D-1+\beta} \frac{1}{\lambda_{dk}^{\frac{1}{2}}} \exp \left[\frac{\lambda_{dk} \mu_{dk}^2}{2} \right]
\end{aligned}$$

1.3 Sampling β

According to reference [?, ?], the prior distribution of the number of variable clusters K can be written as

$$p(K|\beta, D) \propto c_D(K) D! \beta^K \frac{\Gamma(\beta)}{\Gamma(\beta + D)} \quad (15)$$

and $c_D(K) = p(K|\beta = 1, D)$, can be computed by using recurrence formula for Stirling numbers. For $\beta > 0$, the Gamma function in (15) can be computed as

$$\frac{\Gamma(\beta)}{\Gamma(\beta + D)} = \frac{(\beta + D) \mathcal{B}(\beta + 1, D)}{\beta \Gamma(D)} \quad (16)$$

where $\mathcal{B}(\beta + 1, D)$ is a Beta distribution with parameter $\beta + 1, D$.

378 We put a Gamma prior on β and compute the posterior distribution of β as follows:
 379

$$\begin{aligned}
 380 \quad p(\beta|\mathbf{X}, K) &\propto p(\beta|K) \\
 381 \quad &\propto p(K|\beta)p(\beta) \\
 382 \quad &\propto \beta^K \frac{(\beta + D)\mathcal{B}(\beta + 1, D)}{\beta\Gamma(D)} p(\beta)
 \end{aligned} \tag{17}$$

383
 384
 385
 386 By introducing an auxiliary variable u , we can rewrite this in the form:
 387

$$\begin{aligned}
 388 \quad p(\beta|\mathbf{X}, K) &\propto \beta^K \frac{(\beta + D)\mathcal{B}(\beta + 1, D)}{\beta\Gamma(D)} p(\beta) \\
 389 \quad &\propto p(\beta)\beta^{K-1}(\beta + D) \int_0^1 u^\beta(1-u)^{D-1} du
 \end{aligned} \tag{18}$$

390
 391
 392
 393
 394 Hence we can write the joint distribution of β and u as
 395

$$396 \quad p(\beta, u|x, K) \propto p(\beta)\beta^{K-1}(\beta + D)u^\beta(1-u)^{D-1} \quad (\beta \in (0, +\infty], u \in (0, 1)) \tag{19}$$

397
 398
 399 Given the Gamma prior $\mathcal{G}(a_\beta, b_\beta)$ for β , we calculate the conditional distribution of β and u as
 400 follows:
 401

$$\begin{aligned}
 402 \quad p(u|\beta, K) &\propto u^\beta(1-u)^{D-1} \quad (u \in (0, 1)) \\
 403 \quad p(\beta|K, u) &\propto \beta^{K-1}(\beta + D)u^\beta(1-u)^{D-1}\mathcal{G}(a_\beta, b_\beta) \\
 404 \quad &\propto \beta^{K-1}(\beta + D)u^\beta(1-u)^{D-1} \frac{b_\beta^{a_\beta}}{\Gamma(a_\beta)} \beta^{a_\beta-1} \exp(-\beta b_\beta) \\
 405 \quad &\propto \beta^{a_\beta+K-2}(\beta + D)e^{-\beta(b_\beta - \log(u))} \\
 406 \quad &\propto \beta^{a_\beta+K-1}e^{-\beta(b_\beta - \log(u))} + D\beta^{a_\beta+K-2}e^{-\beta(b_\beta - \log(u))} \\
 407 \quad &\propto \frac{\Gamma(a_\beta + K)}{[b_\beta - \log(u)]^{a_\beta+K}} \frac{[b_\beta - \log(u)]^{a_\beta+K}}{\Gamma(a_\beta + K)} \beta^{a_\beta+K-1} e^{-\beta(b_\beta - \log(u))} \\
 408 \quad &+ D \frac{\Gamma(a_\beta + K - 1)}{[b_\beta - \log(u)]^{a_\beta+K-1}} \frac{[b_\beta - \log(u)]^{a_\beta+K-1}}{\Gamma(a_\beta + K - 1)} \beta^{a_\beta+K-2} e^{-\beta(b_\beta - \log(u))} \\
 409 \quad &\propto \frac{\Gamma(a_\beta + K)}{[b_\beta - \log(u)]^{a_\beta+K}} \mathcal{G}(a_\beta + K, b_\beta - \log(u)) \\
 410 \quad &+ D \frac{\Gamma(a_\beta + K - 1)}{[b_\beta - \log(u)]^{a_\beta+K-1}} \mathcal{G}(a_\beta + K - 1, b_\beta - \log(u))
 \end{aligned} \tag{20}$$

411
 412
 413
 414
 415
 416
 417
 418
 419
 420 According to the conditional posteriors derived above, we can sample β and u as follows:
 421

$$\begin{aligned}
 422 \quad (u|\beta, K) &\sim \mathcal{B}(\beta + 1, D) \\
 423 \quad (\beta|u, K) &\sim \pi_K \mathcal{G}(a_\beta + K, b_\beta - \log(u)) + (1 - \pi_K) \mathcal{G}(a_\beta + K - 1, b_\beta - \log(u))
 \end{aligned} \tag{21}$$

424
 425
 426 where the posterior distribution of β is a mixture of two gamma densities with the weights π_K
 427 defined by
 428

$$429 \quad \frac{\pi_K}{(1 - \pi_K)} = \frac{(a_\beta + K - 1)}{D(b_\beta - \log(u))} \tag{22}$$

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

1.4 Sampling the factor matrix F

Given our model is represented as follows:

$$\begin{aligned} p(\mathbf{f}_n) &= \mathcal{N}(\mathbf{f}_n | \mathbf{0}, \mathbf{I}) \\ p(\mathbf{x}_n) &= \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu} + (\mathbf{G} \odot \mathbf{Z})\mathbf{f}_n, \boldsymbol{\Psi}) \end{aligned} \quad (23)$$

We can calculate the joint distribution of \mathbf{f}_n and \mathbf{x}_n as follows:

$$p\left(\begin{bmatrix} \mathbf{f}_n \\ \mathbf{x}_n \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f}_n \\ \mathbf{x}_n \end{bmatrix} \mid \begin{bmatrix} \mathbf{0} \\ \boldsymbol{\mu} \end{bmatrix}, \begin{bmatrix} \mathbf{I} & (\mathbf{G} \odot \mathbf{Z})^T \\ (\mathbf{G} \odot \mathbf{Z}) & (\mathbf{G} \odot \mathbf{Z})(\mathbf{G} \odot \mathbf{Z})^T + \boldsymbol{\Psi} \end{bmatrix}\right) \quad (24)$$

Apply the Gaussian conditioning formulas to the joint distribution derived above, we can get the conditional distribution of the factor matrix F as follows:

$$\begin{aligned} p(\mathbf{f}_n | \mathbf{x}_n) &= \mathcal{N}(\mathbf{f}_n | \boldsymbol{\mu}_f, \boldsymbol{\Sigma}_f) \\ \boldsymbol{\Sigma}_f &= \mathbf{I} - (\mathbf{G} \odot \mathbf{Z})^T ((\mathbf{G} \odot \mathbf{Z})(\mathbf{G} \odot \mathbf{Z})^T + \boldsymbol{\Psi})^{-1} (\mathbf{G} \odot \mathbf{Z}) \\ \boldsymbol{\mu}_f &= (\mathbf{G} \odot \mathbf{Z})^T ((\mathbf{G} \odot \mathbf{Z})(\mathbf{G} \odot \mathbf{Z})^T + \boldsymbol{\Psi})^{-1} \mathbf{x}_n \end{aligned} \quad (25)$$

1.5 Sampling σ_ϵ^2

$\boldsymbol{\Psi} = [\sigma_\epsilon^2 \mathbf{I}]$ is a diagonal matrix. We place a inverse Gamma prior on σ_ϵ^2 as:

$$\begin{aligned} p(\sigma_\epsilon^2 | a_\epsilon, b_\epsilon) &= \mathcal{IG}(a_\epsilon, b_\epsilon) \\ &= \frac{b_\epsilon^{a_\epsilon}}{\Gamma(a_\epsilon)} (\sigma_\epsilon^2)^{-a_\epsilon-1} \exp\left[-\frac{b_\epsilon}{\sigma_\epsilon^2}\right] \end{aligned} \quad (26)$$

We can get that the posterior is also inverse Gamma distributed:

$$\begin{aligned} p(\sigma_\epsilon^2 | x_{dn}, g_{dk}, z_{dk}, f_{kn}, a_\epsilon, b_\epsilon) &= \prod_{n=1}^N \prod_{d=1}^D \prod_{k=1}^K \frac{p(x_{dn} | g_{dk}, z_{dk}, f_{kn}, \sigma_\epsilon^2) p(\sigma_\epsilon^2 | a_\epsilon, b_\epsilon)}{p(x_{dn} | g_{dk}, z_{dk}, f_{kn}, a_\epsilon, b_\epsilon)} \\ &\propto \prod_{n=1}^N \prod_{d=1}^D \prod_{k=1}^K p(x_{dn} | g_{dk}, z_{dk}, f_{kn}, \sigma_\epsilon^2) p(\sigma_\epsilon^2 | a_\epsilon, b_\epsilon) \\ &= \prod_{n=1}^N \prod_{d=1}^D \prod_{k=1}^K \left(\frac{1}{(2\pi\sigma_\epsilon^2)} \exp\left[-\frac{(x_{dn} - g_{dk}f_{kn})^2}{2\sigma_\epsilon^2}\right] \right)^{z_{dk}} \frac{b_\epsilon^{a_\epsilon}}{\Gamma(a_\epsilon)} (\sigma_\epsilon^2)^{-a_\epsilon-1} \exp\left[-\frac{b_\epsilon}{\sigma_\epsilon^2}\right] \\ &= (\sigma_\epsilon^2)^{-a_\epsilon - \frac{ND}{2} - 1} \exp\left[-\frac{b_\epsilon + \frac{\sum_k^K (\sum_{d=1}^D \sum_{n=1}^N (x_{dn} - g_{dk}f_{kn})^2)^{z_{dk}}}{2}}{\sigma_\epsilon^2}\right] \\ &= \mathcal{IG}\left(a_\epsilon + \frac{ND}{2}, b_\epsilon + \frac{\sum_k^K (\sum_{d=1}^D \sum_{n=1}^N (x_{dn} - g_{dk}f_{kn})^2)^{z_{dk}}}{2}\right) \\ &= \mathcal{IG}\left(a_\epsilon + \frac{ND}{2}, b_\epsilon + \frac{\text{tr}([\mathbf{X} - (\mathbf{G} \odot \mathbf{Z})\mathbf{F}]^T [\mathbf{X} - (\mathbf{G} \odot \mathbf{Z})\mathbf{F}])}{2}\right) \end{aligned} \quad (27)$$

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

1.6 Sampling σ_g^2

σ_g^2 is given a inverse-Gamma prior $\mathcal{IG}(a_g, b_g)$. The local graphical model of the matrix \mathbf{G} is depicted in Figure(3).

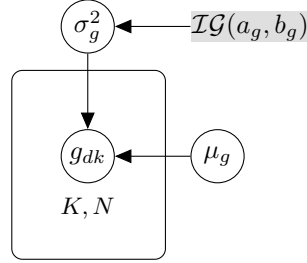


Figure 3: The Local Graphical Model for The Factor Loading Matrix

The posterior distribution of σ_g^2 can be calculated as:

$$\begin{aligned}
 p(\sigma_g^2 | \mathbf{G}, \mu_g, a_g, b_g) &= \prod_{k=1}^K \prod_{d=1}^D p(\sigma_g^2 | g_{dk}, \mu_g, a_g, b_g) \\
 &= \prod_{k=1}^K \prod_{d=1}^D \frac{p(g_{dk} | \mu_g, \sigma_g^2) p(\sigma_g^2 | a_g, b_g)}{p(g_{dk} | \mu_g, a_g, b_g)} \\
 &\propto \prod_{k=1}^K \prod_{d=1}^D p(g_{dk} | \mu_g, \sigma_g^2) p(\sigma_g^2 | a_g, b_g) \\
 &\propto \prod_{k=1}^K \prod_{d=1}^D \frac{1}{(2\pi\sigma_g^2)} \exp\left[-\frac{(g_{dk} - \mu_g)^2}{2\sigma_g^2}\right] \frac{b_g^{a_g}}{\Gamma(a_g)} (\sigma_g^2)^{-a_g-1} \exp\left[-\frac{b_g}{\sigma_g^2}\right] \\
 &= (\sigma_g^2)^{-a_g - \frac{DK}{2} - 1} \exp\left[-\frac{b_g + \frac{\sum_{d=1, k=1}^{D, K} (g_{dk} - \mu_g)^2}{2}}{\sigma_g^2}\right] \\
 &= \mathcal{IG}\left(a_g + \frac{DK}{2}, b_g + \frac{\sum_{d=1, k=1}^{D, K} (g_{dk} - \mu_g)^2}{2}\right)
 \end{aligned}$$

2 The Nonparametric Bayesian Mixture of Sparse Factor Analysers

Consider observed data $\mathbf{X} \in \mathbb{R}^{D \times N}$ where we have D measured variables for each data sample \mathbf{x}_n . Our goal is to cluster the samples into the different groups parameterised by group-specific covariance structure and mean, and for each group of samples, to partition the variables into the clusters where the clustered variables can be interpreted by common latent factors. Our generative model can be represented as:

$$\begin{aligned}
x_{dn} &\sim \sum_k \pi_k \mathcal{N} \left(\mu_d^{(k)} + \sum_j (g_{dj}^{(k)} \circ z_{dj}^{(k)}) f_{jn}, \sigma_d^{(k)2} \right) \\
c_n &\sim \text{Multi}(\pi_1, \dots, \pi_K) \\
\pi_k &\sim \text{GEM}(\alpha) \\
\boldsymbol{\mu}^{(k)} &\sim \mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0) \\
g_{dj}^{(k)} &\sim \mathcal{N}(0, \sigma_g^{(k)2}) \\
c_d^{(k)} &\sim \text{Multi}(\lambda_1^{(k)}, \dots, \lambda_J^{(k)}) \\
\lambda_j^{(k)} &\sim \text{GEM}(\beta) \\
z_{dj}^{(k)} &= \mathbb{1}(c_d = j) \\
f_{jn} &\sim \mathcal{N}(0, I)
\end{aligned} \tag{28}$$

where the samples are drawn from one of the mixture components indicated by variable c_n with the mixing proportions π_k as defined in classical Dirichlet process mixture model. Through factor modelling of each mixture component independently, we can reconstruct the covariances by the sampled factor loading matrices. In our model for the grouped samples linked to a particular mixture component, $c_d^{(k)}$, used to indicate the clustering of variables for samples associated with group k , is drawn from a Dirichlet Process and $g_{dj}^{(k)}$ is used to model the strength of association between variable d and factor j . By introducing indicator variable $z_{dj}^{(k)} = \mathbb{1}(c_d = j)$, the factor loading matrix can be written as $\mathbf{W}^{(k)} = \mathbf{G}^{(k)} \circ \mathbf{Z}^{(k)}$. We use a group-specific parameter $\sigma_d^{(k)}$ to model the idiosyncratic noise. The covariance associated with this component can be modelled as $\Sigma^{(k)} = \mathbf{W}^{(k)} \mathbf{W}^{(k)T} + \sigma_d^{(k)2} I_D$.

2.1 Hyperparameter setting

We put Inverse Gamma priors on $\sigma_g^{(k)2}$ and $\sigma_d^{(k)2}$ independently for each component k and put Gamma priors on α and β .

$$\begin{aligned}
\sigma_g^{(k)2} &\sim \text{IG}(a_g, b_g) \\
\sigma_d^{(k)2} &\sim \text{IG}(a_d, b_d) \\
\alpha &\sim \mathcal{G}(a_\alpha, b_\alpha) \\
\beta &\sim \mathcal{G}(a_\beta, b_\beta)
\end{aligned} \tag{29}$$

3 Inference

We use a Gibbs sampler to exploit the posterior distribution over the model parameters $\boldsymbol{\mu}^{(k)}, g_{dj}^{(k)}, c_d^{(k)}, \mathbf{f}_n, c_n$ and the hyperparameters $\alpha, \beta, \sigma_d^{(k)}, \sigma_g^{(k)}$. The posterior distributions used in Gibbs sampling are summarized here.

Sampling the real-valued matrix $\mathbf{G}^{(j)}$:

As we place a normal distributed prior $\mathcal{N}(0, \sigma_g^{(k)2})$ on the entries $g_{dj}^{(k)}$ of $\mathbf{G}^{(k)}$ independently, we can sample $g_{dj}^{(k)}$ from its posterior as follows:

$$p(g_{dj}^{(k)} | \{\mathbf{x}_n, \mathbf{f}_n\}_{c_n=k}, \sigma_g^{(k)2}, \boldsymbol{\mu}^{(k)}) \sim \mathcal{N}(\mu_{d,j}, \lambda_{d,j}^{-1}) \quad (30)$$

where

$$\lambda_{d,j} = \begin{cases} \frac{\sum_{c_n=k} f_{jn}^2}{\sigma_d^{(k)2}} + \frac{1}{\sigma_g^{(k)2}} & \text{if } c_d = j \\ \frac{1}{\sigma_g^{(k)2}} & \text{otherwise} \end{cases} \quad (31)$$

$$\mu_{d,j} = \begin{cases} \left(\frac{\sum_{c_n=k} f_{jn}^2}{\sigma_d^{(k)2}} + \frac{1}{\sigma_g^{(k)2}} \right)^{-1} \left(\frac{\sum_{c_n=k} (x_{dn} - \mu_d^{(k)}) f_{jn}}{\sigma_d^{(k)2}} \right) & \text{if } c_d = j \\ 0 & \text{otherwise} \end{cases}$$

where $\mu_d^{(k)}$ denote the element d of the mean of mixture component k .

We place an inverse-Gamma prior $\mathcal{IG}(a_g, b_g)$ on the hyperparameter $\sigma_g^{(k)2}$ and its posterior can be calculated as:

$$p(\sigma_g^{(k)2} | \mathbf{G}^{(k)}, a_g, b_g) = \mathcal{IG} \left(a_g + \frac{DJ^{(k)}}{2}, b_g + \frac{\text{tr}(\mathbf{G}^{(k)T} \mathbf{G}^{(k)})}{2} \right)$$

Sampling $\mathbf{Z}^{(j)}$:

We sample the entries $\{z_{dk}^{(j)}\}$ of matrix $\mathbf{Z}^{(j)}$ from its posterior for each mixture component independently:

$$\begin{aligned} p(z_{dk}^{(j)} = 1 | \mathbf{x}_{(c_j)}, g_{dk}^{(j)}, \mathbf{z}_{-dk}^{(j)}, f_{kn}, \beta, \sigma_d^{(j)2}) \\ \propto \prod_{n \in c_j} p(x_{dn} | g_{dk}^{(j)}, z_{dk}^{(j)} = 1, f_{kn}, \beta, \sigma_d^{(j)2}) p(z_{dk}^{(j)} = 1 | \mathbf{z}_{-dk}^{(j)}, \beta) \\ = \frac{1}{(2\pi\sigma_d^{(j)2})^{|c_j|/2}} \exp \left[-\frac{\sum_{n \in c_j} (x_{dn} - g_{dk}^{(j)} f_{kn})^2}{2\sigma_d^{(j)2}} \right] \frac{m_{-dk}^{(j)}}{D - 1 + \beta} \end{aligned} \quad (32)$$

where \mathbf{z}_{-dk} denotes all the indicators except $\mathbf{z}_d = [z_{d1}, \dots, z_{dK}]$, and $m_{-dk}^{(j)}$ denotes the number of variables associated with the latent factor k in component j . **Sampling \mathbf{F} :**

For sampling of the factor matrix \mathbf{F} , we sample the columns of \mathbf{F} independently from its posterior $p(\mathbf{f}_n | \mathbf{x}_n, \mathbf{G}^{(c_n)}, \mathbf{Z}^{(c_n)}, \boldsymbol{\mu}^{(c_n)})$

$$\begin{aligned} p(\mathbf{f}_n | \mathbf{x}_n) &= \mathcal{N}(\mathbf{f}_n | \boldsymbol{\mu}_f^{(c_n)}, \boldsymbol{\Sigma}_f^{(c_n)}) \\ \boldsymbol{\Sigma}_f^{(c_n)} &= \mathbf{I} - (\mathbf{G}^{(c_n)} \odot \mathbf{Z}^{(c_n)})^T \left((\mathbf{G}^{(c_n)} \odot \mathbf{Z}^{(c_n)}) (\mathbf{G}^{(c_n)} \odot \mathbf{Z}^{(c_n)})^T + \sigma_d^{(c_n)2} \mathbf{I} \right)^{-1} (\mathbf{G}^{(c_n)} \odot \mathbf{Z}^{(c_n)}) \\ \boldsymbol{\mu}_f^{(c_n)} &= (\mathbf{G}^{(c_n)} \odot \mathbf{Z}^{(c_n)})^T \left((\mathbf{G}^{(c_n)} \odot \mathbf{Z}^{(c_n)}) (\mathbf{G}^{(c_n)} \odot \mathbf{Z}^{(c_n)})^T + \sigma_d^{(c_n)2} \mathbf{I} \right)^{-1} (\mathbf{x}_n - \boldsymbol{\mu}^{(c_n)}) \end{aligned} \quad (33)$$

Sampling the idiosyncratic noise term: We place an inverse-Gamma prior $\mathcal{IG}(a_d, b_d)$ over $\sigma_d^{(k)2}$ and the posterior is also inverse-Gamma distributed:

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

$$p(\sigma_d^{(k)2} | \cdot) \propto \mathcal{IG} \left(a_d + \frac{N^{(k)}D}{2}, b_d + \frac{\text{tr}(\mathbf{E}^T \mathbf{E})}{2} \right)$$

where $\mathbf{E} = \sum_{c_n=k} (\mathbf{x}_n - \boldsymbol{\mu}^{(k)} - (G^{(k)} \odot Z^{(k)}) \mathbf{f}_n)$.

Sampling $\{c_n\}$: We cluster each data point according to the posterior of the indicator variable c_n that can be calculated as

$$\begin{aligned} & p(c_n = j | \mathbf{x}_n, \mathbf{c}_{-n}, \alpha, \boldsymbol{\mu}^{(j)}, \mathbf{A}^{(j)}, \boldsymbol{\Psi}^{(j)}) \\ & \propto p(c_n = j | \mathbf{c}_{-n}, \alpha) p(\mathbf{x}_n | c_n = j, \boldsymbol{\mu}^{(j)}, \mathbf{A}^{(j)}, \boldsymbol{\Psi}^{(j)}, \mathbf{c}_{-n}) \\ & \propto \frac{N_{-n}^{(j)} |\mathbf{A}^{(j)} \mathbf{A}^{(j)T} + \boldsymbol{\Psi}^{(j)}|^{-\frac{1}{2}}}{N - 1 + \alpha} \\ & \quad \cdot \exp \left(-\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}^{(j)})^T (\mathbf{A}^{(j)} \mathbf{A}^{(j)T} + \boldsymbol{\Psi}^{(j)})^{-1} (\mathbf{x}_n - \boldsymbol{\mu}^{(j)}) \right) \end{aligned} \quad (34)$$

where $|\cdot|$ denotes the determinant of its matrix argument, $N_{-n}^{(j)}$ denotes the number of data samples excluding \mathbf{x}_n that are associated with component j , and \mathbf{c}_{-n} denotes all the indicators for data samples except c_n .

Sampling the idiosyncratic noise term: Since we placed an inverse Gamma prior $\mathcal{IG}(a_\epsilon, b_\epsilon)$ over $\sigma_d^{(j)2}$, the posterior is also inverse Gamma distributed:

$$p(\sigma_d^{(j)2} | -) \propto \mathcal{IG} \left(a_\epsilon + \frac{N^{(j)}D}{2}, b_\epsilon + \frac{\text{tr}(\Gamma^T \Gamma)}{2} \right) \quad (35)$$

where we use $-$ to denote the rest of the model parameters, and $\Gamma = \hat{\mathbf{x}}_{(c_j)} - (G^{(j)} \odot Z^{(j)}) \mathbf{f}_{(c_j)}$.

Sampling hyperparameters: We sample the concentration parameter α from a joint distribution by introducing an auxiliary variable u as

$$\begin{aligned} & p(\alpha | u, J) \sim \gamma_J \mathcal{G}(a_\alpha + J, b_\alpha - \log(u)) \\ & \quad + (1 - \gamma_J) \mathcal{G}(a_\alpha + J - 1, b_\alpha - \log(u)) \\ & p(u | \alpha, N) \sim \text{Beta}(\alpha + 1, N) \end{aligned} \quad (36)$$

with the weights γ_J defined by

$$\frac{\gamma_J}{(1 - \gamma_J)} = \frac{a_\alpha + J - 1}{N(b_\alpha - \log(u))} \quad (37)$$

where N denotes the number of data samples and \log denotes the natural logarithm of its argument.

We use the same sampling procedure to sample the second-layer DP concentration parameter β conditioned upon the maximum number of latent factors in the mixture distribution K and data dimension D .