# Collapsed Variational Inference for Nonparametric Bayesian Group Factor Analysis

**Sikun Yang,   Heinz Koeppl**
Department of Electrical Engineering and Information Technology
Technische Universität Darmstadt
64283 Darmstadt, Germany
{sikun.yang, heinz.koeppl}@bcs.tu-darmstadt.de

## Abstract

Group factor analysis (GFA) methods have been widely used to infer the common structure and the group-specific signals from multiple related datasets in various fields including systems biology and neuroimaging. To date, most available GFA models require Gibbs sampling or slice sampling to perform inference, which prevents the practical application of GFA to large-scale data. In this paper we present an efficient collapsed variational inference (CVI) algorithm for the nonparametric Bayesian group factor analysis (NGFA) model built upon an hierarchical beta Bernoulli process. Our CVI algorithm proceeds by marginalizing out the group-specific beta process parameters, and then approximating the true posterior in the collapsed space using mean field methods. Experimental results on both synthetic and real-world data demonstrate the effectiveness of our CVI algorithm for the NGFA compared with state-of-the-art GFA methods.

## 1   Introduction

Factor analysis (FA) is a powerful tool widely used to infer low-dimensional structure in multivariate data. More specifically, FA models attempt to represent a data matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$ by the product of two matrices plus residual noise as

$$\mathbf{X} = \mathbf{FG} + \mathbf{E},$$

where $\mathbf{F} \in \mathbb{R}^{N \times K}$ denotes the factor score matrix, and $\mathbf{G} \in \mathbb{R}^{K \times D}$ denotes the factor loading matrix; $\mathbf{E} \in \mathbb{R}^{N \times D}$ is the residual noise matrix. For high-dimensional data, FA models imposing sparsity-inducing priors (West, 2003; Rai and Daume III, 2008; Paisley et al., 2009; Knowles et al., 2011) or regularizations (Zou

et al., 2006; Witten et al., 2009) over the inferred loading matrices are developed to improve interpretability of the inferred low-dimensional structure. For example, in gene expression analysis, a factor loading matrix characterizing the connections between transcription factors and regulated genes are expected to be sparse (Carvalho et al., 2008).

In many real-world applications, we often deal with multiple related datasets – each comprising a group of variables – that need to be factorized in a common subspace. For instance, latent Dirichlet allocation (Blei et al., 2003) and Poisson factor analysis models (Zhou et al., 2015) have been developed to learn the shared latent topics among multiple related documents. Recently, GFA models (Virtanen et al., 2012; Bunte et al., 2016) using the automatic relevance determination (ARD) prior have been proposed for drug sensitivity prediction and functional neuroimaging. However, the modeling flexibility achieved by these GFA models comes at a price as their inference usually requires Markov chain Monte Carlo (MCMC) to perform posterior computation, which makes them to scale poorly for large-scale GFA problems. Alternatively, variational Bayesian inference has been shown to be efficient for large-scale data by making an independence assumption among latent variables and parameters (Wainwright et al., 2008). However, this strong assumption may lead to very inaccurate results in practical applications, especially for GFA problems where latent variables might be tightly coupled.

Motivated by this limitation, we propose a computationally efficient collapsed variational inference algorithm for the nonparametric Bayesian group factor analysis model. Our NGFA model is built upon the hierarchical beta process (HBP) (Thibaux et al., 2007). We note that the HBP has been investigated in (Chen et al., 2011; Gupta et al., 2012a,b) for joint modeling of multiple data matrices utilizing MCMC, but again showed poor scalability and slow convergence. For nonparametric Bayesian models, such as HDP topic model (Teh et al., 2007) and HDP hid-
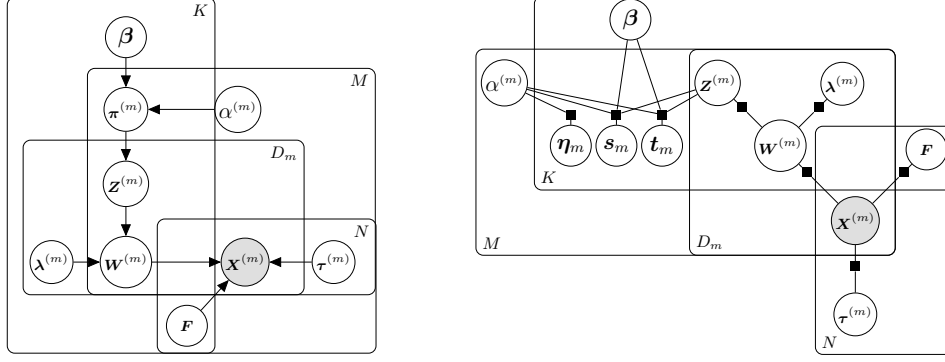
Figure 1: Left: The graphical representation of the proposed model. Right: Factor graph of the model with auxiliary variables.

den Markov models (Fox et al., 2011), collapsed Gibbs sampling (CGS) are typically employed to perform posterior computation because CGS rapidly convergences onto the true posterior. However, it remains challenging to assess the convergence of CGS algorithms for practical use. To address this issue, collapsed variational inference algorithms (Teh et al., 2006, 2008; Foulds et al., 2013) are developed for topic models by integrating out model parameters, and then applying the mean field approximation to the latent variables. Recently, collapsed variational inference algorithms have been developed for hidden Markov models (Wang et al., 2013), nonparametric relational models (Ishiguro et al., 2017) and Markov jump processes (Zhang et al., 2017) with encouraging results. In this paper, we aim to develop a collapsed variational inference algorithm for the nonparametric Bayesian group factor analysis model.

We make the following contributions:

- We tackle the group factor analysis problems using a Bayesian nonparametric method based on the hierarchical beta Bernoulli process. The total number of factors is automatically learned from data. Specifically, the NGFA model induces both group-wise and element-wise structured sparsity effectively compared to state-of-the-art GFA methods (see Section 4.1).
- An efficient collapsed variational inference algorithm is proposed to infer the NGFA model.
- We apply the developed method to real world multiple related dataset, with encouraging results (see Section 4.2; 4.3).

The paper is organized as follows. In Section 2, we describe the nonparametric Bayesian group factor analysis model. Our collapsed variational inference algorithm for the NGFA is introduced in Section 3. Experimental results are presented in Section 4. Finally, conclusions and

possible directions for future research are discussed in Section 5.

## 2 Nonparametric Bayesian Group Factor Analysis

Given multiple related data matrices $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \ldots, \mathbf{X}^{(M)}$, each with $N$ samples, i.e., $\mathbf{X}^{(m)} \in \mathbb{R}^{N \times D_m}$, our goal is to factorize each dataset $\mathbf{X}^{(m)}$ into the product of a common factor matrix $\mathbf{F} = [\mathbf{f}_1, \ldots, \mathbf{f}_K]$ of size $N \times K$, and a group-specific factor loading matrix $\mathbf{G}^{(m)}$ of size $K \times D_m$ as

$$\mathbf{X}^{(m)} = \mathbf{F}\mathbf{G}^{(m)} + \mathbf{E}^{(m)}, \tag{1}$$

where $\mathbf{E}^{(m)} = [\mathbf{e}_1^{(m)}, \ldots, \mathbf{e}_{D_m}^{(m)}]$ is assumed to be Gaussian noise for the $m$-th dataset or group. We impose independent normal priors over $\mathbf{e}_d^{(m)} \in \mathbb{R}^N$, i.e., $\mathbf{e}_d^{(m)} \sim \mathcal{N}(0, \text{diag}(\tau_1^{(m)}, \ldots, \tau_N^{(m)}))$, where $\tau_n^{(m)}$ controls the variance of $N$-th sample in the $m$-th group. As commonly used in factor analysis (Rai and Daume III, 2008; Paisley et al., 2009; Knowles et al., 2011), we put a normal prior on each factor $\mathbf{f}_k$, i.e., $\mathbf{f}_k \sim \mathcal{N}(0, \mathbf{I}_N)$, where $\mathbf{I}_N$ is an identity matrix of size $N$. To explicitly capture the sparsity, we model the factor loading matrix $\mathbf{G}^{(m)}$ for each group by the element-wise product of a binary matrix $\mathbf{Z}^{(m)}$ and a real-valued weight matrix $\mathbf{W}^{(m)}$, i.e., $\mathbf{G}^{(m)} = \mathbf{Z}^{(m)} \odot \mathbf{W}^{(m)}$. More specifically, we place a normal prior over each element of $\mathbf{W}^{(m)}$, i.e., $w_{kd}^{(m)} \sim \mathcal{N}(0, (\lambda_{kd}^{(m)})^{-1})$. To allow the number of factors $K$ to be automatically inferred from data, we model each row of $\mathbf{Z}^{(m)}$ as a draw from a group-specific Bernoulli process. As our goal is to factorize multiple related data matrices using a common set of factors, we naturally consider the hierarchical beta process (Thibaux et al., 2007) that allows us to generate a set of latent factors from a *global* beta process $B$, and then allow the generated factors to be shared among all the groups. The usage of

$$p(\mathbf{Z} \mid \boldsymbol{\beta}, \boldsymbol{\alpha}) = \int p(\mathbf{Z} \mid \boldsymbol{\pi}) p(\boldsymbol{\pi} \mid \boldsymbol{\beta}, \boldsymbol{\alpha}) \mathrm{d}\boldsymbol{\pi} = \prod_{m,k} \frac{\Gamma(\alpha^{(m)})}{\Gamma(\alpha^{(m)} + D_m)} \frac{\Gamma(\alpha^{(m)}\beta_k + \hat{n}_{mk})}{\Gamma(\alpha^{(m)}\beta_k)} \frac{\Gamma(\alpha^{(m)}\bar{\beta}_k + \tilde{n}_{mk})}{\Gamma(\alpha^{(m)}\bar{\beta}_k)}$$

Figure 2: The marginal distribution of $\mathbf{Z}$. We define $\hat{n}_{mk} \equiv \sum_d \mathbb{1}(z_{kd}^{(m)} = 1)$ and $\tilde{n}_{mk} \equiv \sum_d \mathbb{1}(z_{kd}^{(m)} = 0)$, where $\mathbb{1}(\cdot)$ is the standard indicator function.

the generated factors in each group is determined by the group-specific beta process $A^{(m)}$. More specifically, the hierarchical beta Bernoulli process is

$$B \equiv \sum_{k=1}^{K} \beta_k \delta_{\mathbf{f}_k}, \qquad A^{(m)} \equiv \sum_{k=1}^{K} \pi_k^{(m)} \delta_{\mathbf{f}_k}, \qquad (2)$$
$$\mathbf{f}_k \sim \mathcal{N}(0, \mathbf{I}_N), \qquad \beta_k \sim \mathrm{Beta}(\kappa_0/K, \kappa_0(K-1)/K),$$
$$\pi_k^{(m)} \sim \mathrm{Beta}(\alpha^{(m)}\beta_k, \alpha^{(m)}\bar{\beta}_k),$$
$$z_{kd}^{(m)} \sim \mathrm{Bern}(\pi_k^{(m)}),$$

where $\bar{\beta}_k \equiv 1 - \beta_k$, and $K$ is a truncation level that is set sufficiently large to ensure a good approximation to the truly infinite model. The concentration parameters of the global beta process and the local group-specific beta process are $\kappa_0$ and $\alpha^{(m)}$, respectively. The total number of factors shared among all groups is determined by $\kappa_0$, and the amount of variability of each $A^{(m)}$ around $B$ is determined by $\alpha^{(m)}$. To improve the flexility of the model, we place gamma priors on $\lambda_{kd}^{(m)}$, $\tau_n^{(m)}$ and $\alpha^{(m)}$, respectively, as $\lambda_{kd}^{(m)} \sim \mathrm{Gam}(g_0, h_0)$, $\tau_n^{(m)} \sim \mathrm{Gam}(e_0, f_0)$, $\alpha^{(m)} \sim \mathrm{Gam}(c_0, d_0)$. The graphical representation of the NGFA model is shown in shown in Fig. 1 (left).

## 3 Collapsed Variational Inference

The main idea of collapsed variational inference is to marginalize out model parameters, and then apply the mean field method to approximate the distribution over latent variables. We note that marginalizing out the parameters induces dependencies among the latent variables. However, each latent variable interacts with the remaining variables only through the sufficient statistics (i.e. the field) in the collapsed space, and the influence of any single variable on the field is small. Hence, the dependency between any two latent variables is weak, suggesting that the mean field assumption is better justified in the collapsed space. In our case, we first marginalize out the group-specific beta process parameters to obtain the marginal distribution over latent variables. We then employ the variational posterior to approximate the distribution of latent variables and the remaining parameters.
**Notation.** When expressing the conditional distribution, we will use the shorthand "–" to denote full conditionals, i.e., all other variables. For the sake of clarity, we use

$\mathbf{X}$ to denote the set of matrices $(\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(M)})$. Similarly, let $\mathbf{Z}$ denote $(\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(M)})$, and $\boldsymbol{\pi}$ denote $(\boldsymbol{\pi}^{(1)}, \dots, \boldsymbol{\pi}^{(M)})$. With slight notational abuse we use generic $p$ to denote probability density and mass functions.

We repeatedly exploit the following three results (Teh et al., 2008) to derive the collapsed variational inference algorithm for the NGFA.

**Result 1.** The geometric expectation of a non-negative random variable $y$ is defined as $\mathsf{G}[y] \equiv \exp(\mathsf{E}[\log(y)])$. If $y$ is gamma distributed, i.e., $p(y \mid a, b) \propto y^{a-1}e^{-by}$, the geometric expectation of $y$ is $\mathsf{G}[y] = \frac{\exp(\Psi(a))}{b}$, where $\Psi(y) = \frac{\partial \log \Gamma(y)}{\partial y}$ is the digamma function. For a beta distributed random variable $y$, i.e., $p(y \mid a, b) \propto y^{a-1}(1-y)^{b-1}$, the geometric expectation of $y$ is $\mathsf{G}[y] = \frac{\exp[\Psi(a)]}{\exp[\Psi(a+b)]}$. If $y_1, \dots, y_K$ are mutually independent, we have, $\mathsf{G}\left[\prod_{k=1}^{K} y_k\right] = \prod_{k=1}^{K} \mathsf{G}[y_k]$.

**Result 2.** According to the central limit theorem, if $y$ is the sum of $N$ independent Bernoulli random variables, i.e., $y = \sum_{i=1}^{N} u_i$, where $u_i \sim \mathrm{Bern}(\xi_i)$, then for large enough $N$, $y$ is well approximated by a Gaussian random variable with mean and variance as

$$\mathsf{E}[y] = \sum_{i=1}^{N} \xi_i, \qquad \mathsf{V}[y] = \sum_{i=1}^{N} \xi_i(1-\xi_i),$$

respectively. Moreover, the expectation of $\log(y)$ can be approximated using the second-order Taylor expansion (Hoef, 2012) as

$$\mathsf{E}[\log(y)] \approx \log(\mathsf{E}[y]) - \frac{\mathsf{V}[y]}{2(\mathsf{E}[y])^2}.$$

**Result 3.** If $l$ is the sum of independent Bernoulli random variables, i.e., $l = \sum_i u_i$, where $u_i \sim \mathrm{Bern}(\xi_i)$, we use $p_+(l)$ to denote the probability of $l$ being positive, i.e.,

$$p_+(l) \equiv p(l > 0) = 1 - \prod_i p(u_i = 0)$$

$$= 1 - \exp\left[\sum_i \log(1 - \xi_i)\right].$$

Accordingly, the expectation and variance conditional on $l > 0$ are defined as $\mathsf{E}_+[l] \equiv \frac{\mathsf{E}[l]}{p_+(l)}$ and $\mathsf{V}_+[l] \equiv \frac{\mathsf{V}[l]}{p_+(l)}$,

respectively. If $y$ is then a Chinese restaurant table (CRT) (Pitman, 2006) distributed random variable, i.e., $p(y \mid a, l) = \frac{\Gamma(a)}{\Gamma(a+l)} \begin{bmatrix} l \\ y \end{bmatrix} a^y$, where $y = 0, 1, \ldots, l$, and $\begin{bmatrix} n \\ m \end{bmatrix}$ denoting the unsigned Stirling number of the first kind, then the expectation of $y$ can be closely approximated using the improved second-order Taylor expansion as

$$
\begin{aligned}
\mathsf{E}[y] \approx \mathsf{G}[a] p_+(l) \Big( &\Psi\big(\mathsf{G}[a] + \mathsf{E}_+[l]\big) \\
&- \Psi(\mathsf{G}[a]) + \frac{\mathsf{V}_+[l]\Psi'(\mathsf{G}[a] + \mathsf{E}_+[l])}{2} \Big),
\end{aligned}
$$

where $\Psi'(y) = \frac{\partial^2 \log \Gamma(y)}{\partial y^2}$ is the trigamma function.

### 3.1 Collapsed representation

First, we describe how to obtain the marginal distribution of latent variables. In the next subsection, we will then describe how to derive the CVI algorithm in the collapsed space.

For the NGFA introduced in the previous section, integrating out $\boldsymbol{\pi}$ yields the marginal distribution of $\mathbf{Z}$ shown in Fig. 2 because beta priors are conjugate to Bernoulli distributions. As the ratios of gamma functions in Fig. 2 give rise to difficulties for updating hyperparameter posteriors, we augment the marginal distribution $\mathbf{Z}$ by introducing three sets of auxiliary variables. More specifically, using the auxiliary variable method (Teh et al., 2007), the first ratio of gamma function can be re-expressed as

$$
\begin{aligned}
&\frac{\Gamma(\alpha^{(m)})}{\Gamma(\alpha^{(m)} + D_m)} \\
&= \frac{1}{\Gamma(D_m)} \int_0^1 \eta_m^{\alpha^{(m)}} (1 - \eta_m)^{D_m - 1} \left(1 + \frac{D_m}{\alpha^{(m)}}\right) \mathrm{d}\eta_m.
\end{aligned} \tag{3}
$$

Via the relation between the gamma function and the Stirling numbers of the first kind (Teh et al., 2007), the second and third ratio of gamma functions can be re-expressed, respectively, as

$$
\frac{\Gamma(\alpha^{(m)}\beta_k + \hat{n}_{mk})}{\Gamma(\alpha^{(m)}\beta_k)} = \sum_{s_{mk}=0}^{\hat{n}_{mk}} \begin{bmatrix} \hat{n}_{mk} \\ s_{mk} \end{bmatrix} (\alpha^{(m)}\beta_k)^{s_{mk}}, \tag{4}
$$

$$
\frac{\Gamma(\alpha^{(m)}\bar{\beta}_k + \tilde{n}_{mk})}{\Gamma(\alpha^{(m)})} = \sum_{t_{mk}=0}^{\tilde{n}_{mk}} \begin{bmatrix} \tilde{n}_{mk} \\ t_{mk} \end{bmatrix} (\alpha^{(m)}\bar{\beta}_k)^{t_{mk}}. \tag{5}
$$

Substituting (Eqs. 3; 4; 5) into Fig. 2, we immediately obtain the joint distribution of the latent and auxiliary variables as

$$
\begin{aligned}
p(\mathbf{Z}, \mathbf{s}, \mathbf{t}, \boldsymbol{\eta} \mid \boldsymbol{\beta}, \boldsymbol{\alpha}) \propto \prod_{m,k} &\eta_m^{\alpha^{(m)}-1}(1-\eta_m)^{D_m-1} \tag{6} \\
&\times \begin{bmatrix} \hat{n}_{mk} \\ s_{mk} \end{bmatrix} (\alpha^{(m)}\beta_k)^{s_{mk}} \begin{bmatrix} \tilde{n}_{mk} \\ t_{mk} \end{bmatrix} (\alpha^{(m)}\bar{\beta}_k)^{t_{mk}}.
\end{aligned}
$$

The factor graph of the expanded system with auxiliary variables is shown in Fig. 1 (right). The conditional distribution of a single latent variable $z_{kd}^{(m)}$ can be derived using the marginal distribution of $\mathbf{Z}$ and the likelihood function according to Eq. 1 as

$$
p(z_{kd}^{(m)} = 1 \mid -) \propto \exp\left[\log(\alpha^{(m)}\beta_k + \hat{n}_{km}^{\neg d})\right] \tag{7}
$$

$$
\times \exp\left[-\frac{1}{2}\sum_n \tau_n^{(m)} \left((w_{kd}^{(m)})^2 f_{nk}^2 - 2w_{kd}^{(m)}\, \tilde{x}_{nd}^{(m)\,\neg k}\right)\right],
$$

where $(\tilde{x}_{nd}^{(m)})^{\neg k} \equiv (x_{nd}^{(m)} - \sum_{j \neq k} z_{jd}^{(m)} w_{jd}^{(m)} f_{nj})$, and $\hat{n}_{km}^{\neg d} \equiv \sum_{d' \neq d} \mathbb{1}(z_{kd'}^{(m)} = 1)$.

### 3.2 Variational approximation

Next, we shall introduce the variational approximation for our expanded system. For the sake of simplicity, the remaining parameters $(\mathbf{W}, \mathbf{F}, \boldsymbol{\beta}, \boldsymbol{\lambda}, \boldsymbol{\tau}, \boldsymbol{\alpha})$ is denoted by $\boldsymbol{\theta}$. Formally, the variational posterior over the augmented variables system is assumed to be of the form

$$
q(\mathbf{Z}, \boldsymbol{\theta}, \mathbf{s}, \mathbf{t}, \boldsymbol{\eta}) = q(\boldsymbol{\theta})q(\mathbf{s}, \mathbf{t}, \boldsymbol{\eta} \mid \mathbf{Z})q(\mathbf{Z}),
$$

where $q(\boldsymbol{\theta}) \equiv q(\mathbf{W})q(\mathbf{F})q(\boldsymbol{\beta})q(\boldsymbol{\lambda})q(\boldsymbol{\tau})q(\boldsymbol{\alpha})$. Note that the true posterior $p(\mathbf{s}, \mathbf{t}, \boldsymbol{\eta} \mid \mathbf{Z})$ is used in our variational update subsequently.

**Evidence Lower Bound (ELBO):** The log marginal likelihood of data is lower bounded as

$$
\begin{aligned}
\log p(\mathbf{X} \mid \kappa_0) &\geq \mathcal{L}(q(\boldsymbol{\theta})q(\mathbf{s}, \mathbf{t}, \boldsymbol{\eta} \mid \mathbf{Z})q(\mathbf{Z})) \tag{9} \\
&= \mathsf{E}_{q(\boldsymbol{\theta}, \mathbf{Z})}\left[\log p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\theta} \mid \kappa_0) - q(\boldsymbol{\theta}, \mathbf{Z})\right].
\end{aligned}
$$

See the supplementary material (A.3) for details.

To maximize the ELBO in Eq. 9 with respect to the variational parameters, we can take the gradients of the ELBO w.r.t. each parameter, and set it equal to zero. Then, our CVI algorithm proceeds by updating the variational parameters in a coordinate-wise manner.

**Updating $q(\mathbf{Z})$:** The variational update for each latent variable $z_{kd}^{(m)}$ is

$$
q(z_{kd}^{(m)} = 1) \propto \exp\left(\mathsf{E}_{q(\mathbf{Z}, \boldsymbol{\theta} \setminus z_{kd}^{(m)})}\left[\log p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\theta} \mid \kappa_0)\right]\right)
$$

$$
\propto \exp\left(\mathsf{E}_{q(\mathbf{Z}, \boldsymbol{\theta} \setminus z_{kd}^{(m)})}\left[\log p(z_{kd}^{(m)} = 1 \mid -)\right]\right), \tag{10}
$$

where $(\mathbf{Z}, \boldsymbol{\theta} \setminus z_{kd}^{(m)})$ means all the variables and parameters excluding $z_{kd}^{(m)}$.

Plugging Eq. 7 into Eq. 10, we obtain the variational update for $q(z_{kd}^{(m)} = 1)$ in Fig. 3. The exact computation of the log count in Fig. 3 is too expensive in practice. According to Result 2, we can approximate it as

$$
\begin{aligned}
\mathsf{E}\left[\log\left(\alpha^{(m)}\beta_k + \hat{n}_{mk}^{\neg d}\right)\right] \approx{}& \log\left(\mathsf{G}[\alpha^{(m)}\beta_k] + \mathsf{E}\left[\hat{n}_{mk}^{\neg d}\right]\right) \\
&- \frac{\mathsf{V}\left[\hat{n}_{mk}^{\neg d}\right]}{2\left(\mathsf{G}[\alpha^{(m)}\beta_k] + \mathsf{E}\left[\hat{n}_{mk}^{\neg d}\right]\right)^2},
\end{aligned}
$$

$$q(z_{kd}^{(m)} = 1) \propto \exp\left\{ \mathsf{E}\left[\log\left(\alpha^{(m)}\beta_k + \hat{n}_{mk}^{\neg d}\right)\right] - \frac{1}{2}\sum_n \mathsf{E}\left[\tau_n^{(m)}\right]\left(\mathsf{E}\left[(w_{kd}^{(m)})^2\right]\mathsf{E}\left[f_{nk}^2\right] - 2\mathsf{E}\left[w_{kd}^{(m)}\right]\mathsf{E}\left[f_{nk}\right]\tilde{x}_{nd}^{(m)\;\neg k}\right)\right\} \tag{8}$$

Figure 3: The variational update for each latent variable.

where the mean and variance of $\hat{n}_{mk}^{\neg d}$ are given by

$$\mathsf{E}\left[\hat{n}_{mk}^{\neg d}\right] = \sum_{d' \neq d} q(z_{kd}^{(m)} = 1),$$

$$\mathsf{V}\left[\hat{n}_{mk}^{\neg d}\right] = \sum_{d' \neq d} q(z_{kd}^{(m)} = 1) q(z_{kd}^{(m)} = 0).$$

**Updating auxiliary variables:** Now we explain how to update the auxiliary variables efficiently using Gaussian approximation techniques. The variational posteriors for the auxiliary variables $\boldsymbol{\eta}$ is

$$q(\boldsymbol{\eta} \mid \mathbf{Z}) \propto \prod_m \eta_m^{\mathsf{E}[\alpha^{(m)}]-1}(1 - \eta_m)^{D_m - 1}.$$

As $\boldsymbol{\eta}$ is beta distributed, via the geometric expectation of Result 1, we have

$$\mathsf{E}[\log(\eta_m)] = \log\left[\mathsf{G}(\eta_m)\right] = \Psi(\mathsf{E}[\alpha^{(m)}]) - \Psi(\mathsf{E}[\alpha^{(m)}] + D_m)$$

The variational posteriors for the auxiliary variables $\mathbf{s}$ is

$$q(\mathbf{s} \mid \mathbf{Z}) \propto \prod_{m,k} \begin{bmatrix} \hat{n}_{mk} \\ s_{mk} \end{bmatrix} (\mathsf{G}[\alpha^{(m)}\beta_k])^{s_{mk}}, \tag{11}$$

where the expectation of $\mathbf{s}$ depends on $\mathbf{Z}$ through the count $\hat{n}_{mk}$ that can take many values. Hence, the exact computation of Eq. 11 is too expensive. According to Result 3, we use the improved second-order Taylor expansion to approximate the expectation of $s_{mk}$ as

$$\mathsf{E}[s_{mk}] \approx \mathsf{G}[\alpha^{(m)}\beta_k] p_+(\hat{n}_{mk})\Big(\Psi\big(\mathsf{G}[\alpha^{(m)}\beta_k] + \mathsf{E}_+[\hat{n}_{mk}]\big)$$

$$- \Psi(\mathsf{G}[\alpha^{(m)}\beta_k]) + \frac{\mathsf{V}_+[\hat{n}_{mk}]\Psi'(\mathsf{G}[\alpha^{(m)}\beta_k] + \mathsf{E}_+[\hat{n}_{mk}])}{2}\Big).$$

Likewise, we can derive the variational update for $\mathbf{t}$ in the same manner. Following the exponential family computation (Wainwright et al., 2008), the variational updates for the remaining parameters are obtained via the conjugacy of our model specification. We present these variational updates in the supplementary material (A.4).

## 4 Experiments

In this section, we compare the nonparametric Bayesian group factor analysis using our proposed CVI algorithm with the state-of-the-art GFA models. We evaluate the proposed CVI algorithm on both synthetic data and real-world applications. In all our experiments, we set $\kappa_0 = 1, c_0 = 0.1, d_0 = 0.1, g_0 = 0.1, h_0 = 0.1, e_0 = 0.1, f_0 = 0.1$. Similar results are obtained when instead setting $\kappa_0 = 0.1, \kappa_0 = 10$ in a sensitivity analysis. Code is available at `https://github.com/stephenyang/CVB_NGFA`.

### 4.1 Simulated data

For our evaluations on synthetic data, we adopt the simulation study in (Zhao et al., 2016): we perform two simulations (*Simulation 1* and *Simulation 2*) which include four groups of data with the dimensionality $D_m = 100$ for each group, respectively. The numbers of samples in the four groups are set to $N = \{20, 40, 60, 100\}$, respectively. In *Simulation 1*, we set the number of latent factors $K = 6$, and generate data only with sparse factor loadings. Specifically, the first three factors are specific to $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}$ and $\mathbf{X}^{(3)}$, respectively, and the last three are shared among all groups. In *Simulation 2*, we set $K = 8$ and generate data with both sparse and dense factor loadings. The sparsity pattern is described in Table 1, and also shown in Fig. 7.

|  | Simulation 1 | | | | | | Simulation 2 | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 | 6 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $X^{(1)}$ | s | - | - | s | - | - | s | - | - | - | d | - | - | - |
| $X^{(2)}$ | - | s | - | s | s | s | - | s | - | s | - | d | - | - |
| $X^{(3)}$ | - | - | s | - | s | s | - | - | s | s | - | - | d | - |
| $X^{(4)}$ | - | - | - | - | - | s | - | - | s | - | - | - | - | d |

Table 1: Sparsity pattern of the factor loading matrices in Simulation 1 and 2. "s" represents a sparse column vector; "d" represents a dense column vector; "-" represents no contribution to that group from the factor.

The sparsity of the sparse factor loadings is handled by setting $90\%$ of the entries in each loading column to zero at random, and the nonzero entries in both the sparse and dense factor loadings are generated from a Gaussian distribution $\mathcal{N}(0, 4)$. The latent factors are generated from a standard Gaussian distribution (i.e., zero mean and unit variance). We generate the residual noise i.i.d. from a Gaussian distribution $\mathcal{N}(0, 1)$.

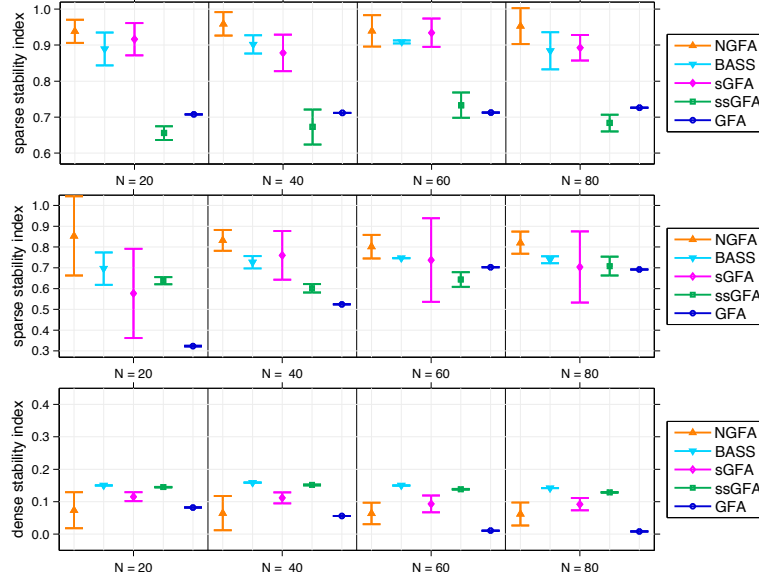We compare the following methods: **(1) GFA:** The

Figure 4: The comparison of stability indices on the inferred matrix of factor loadings for our synthetic data. For SSI, higher is better; for DSI, lower is better. The means and the standard derivations of the stability indices are denoted by the marker and the bar respectively. The SSI comparison of all methods in Simulation 1 is shown in upper rows; The SSI and DSI comparisons in Simulation 2 are shown in middle and bottom rows, respectively.

Bayesian group factor analysis model (Virtanen et al., 2012) with column-wise ARD priors to induce column-wise sparsity on the factor loading matrix. For the GFA model, we used the GFA package with the default parameters setting as set in the code released online. [1] The initial number of factors is set to the true values. The optimization method is L-BFGS with the maximum iterations set to $10^5$. **(2) sGFA:** The extension of the GFA with element-wise ARD priors inducing element-wise sparsity (Bunte et al., 2016). For the sGFA model, the initial number of factors is set to half of the minimum of the sample size and the total number of variables, i.e., $K = \min(N, \sum_m D_m)$. The total number of MCMC iterations is set to $10^5$ with sampling steps set to $10^3$ and thinning steps set to 5. **(3) ssGFA:** The extension of the GFA with the spike-and-slab prior (Bunte et al., 2016), for which we again use the GFA package with the spike-and-slab prior. We set the noise parameters by the informativeNoisePrior function to prevent overfitting. The initial number of factors is set to half of the minimum of the sample size and the total number of variables. The total number of MCMC iterations is set to $10^5$ with sampling steps set to $10^3$ and thinning steps set to 5. **(4) BASS:** The Bayesian group factor analysis with structured sparsity priors (BASS) (Zhao et al., 2016), for which we use the code released in (Zhao et al., 2016). [2] The BASS is initialized using 50 iterations of MCMC and

followed by expectation maximization until convergence, reached when both the number of nonzero loadings do not change for $t$ iterations and the log-likelihood change is less than $1 \times 10^{-5}$ within $t$ iterations. The initial number of factors is set to 10 in *Simulation* 1 and 15 in *Simulation* 2 as described in (Zhao et al., 2016). We perform 20 runs for each method, in particular to evaluate the sensitivity of our inference algorithm to initialization since CVI algorithms are only guaranteed to converge to a local optimum. For all the experiments, we simply set the initial number of factors for our method to be the minimum of the sample size and the dimensionality of each group, and run the model with CVI algorithm until convergence.

To evaluate the performance of the methods on the recovery of sparse and dense factor loadings, we use the sparse and dense stability index defined in (Zhao et al., 2016) to quantify the distance between the true and the inferred factor loading matrices. Given the absolute correlation matrix $\mathbf{C} \in \mathbb{R}^{K_1 \times K_2}$ of the columns of two sparse matrices, the *sparse stability index* (SSI) is calculated as

$$
\text{SSI} = \frac{1}{2K_1} \sum_{r=1}^{K_1} \Big( \max(\mathbf{C}_{r:}) - \frac{\sum_l \mathbb{1}(\mathbf{C}_{rl} > \hat{\mathbf{C}}_{r:})\mathbf{C}_{rl}}{K_2 - 1} \Big)
$$
$$
+ \frac{1}{2K_2} \sum_{l=1}^{K_2} \Big( \max(\mathbf{C}_{:l}) - \frac{\sum_r \mathbb{1}(\mathbf{C}_{rl} > \hat{\mathbf{C}}_{:l})\mathbf{C}_{rl}}{K_1 - 1} \Big),
$$

where $\mathbf{C}_{r:}$ and $\mathbf{C}_{:l}$ denote the $r$-th row and $l$-th column of the matrix $\mathbf{C}$, respectively; $\hat{\mathbf{C}}_{r:}$ and $\hat{\mathbf{C}}_{:l}$ denote the mean of the $r$-th row and $l$-th column of the matrix $\mathbf{C}$,

[1]https://cran.r-project.org/web/packages/GFA/index.html.
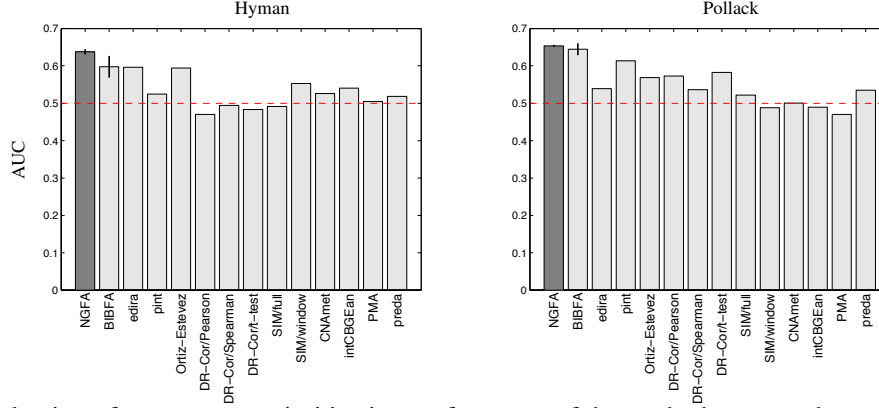[2]https://github.com/judyboon/BASS.

Figure 5: Evaluation of cancer gene prioritization performance of the methods on two data sets: Hyman (left) and Pollack (right). The result is quantified by the area under the ROC curve (AUC). The dashed line indicates the AUC score for a random list (AUC = 0.5). The comparison shows that the NGFA achieves best performance.

respectively. The SSI is invariant to column-scaling and -permutation; larger values indicate better recovery.

The *dense stability index* (DSI) measures the distance between dense matrix columns. Given two dense matrices $\mathbf{M}_1 \in \mathbb{R}^{K_1 \times D}$ and $\mathbf{M}_2 \in \mathbb{R}^{K_2 \times D}$, the DSI is defined as

$$\text{DSI} = \frac{1}{D^2} \text{tr}(\mathbf{M}_1 \mathbf{M}_1^T - \mathbf{M}_2 \mathbf{M}_2^T).$$

The DSI is invariant to orthogonal matrix transformation, column-scaling and -permutation; the lower values indicate better recovery.

Following the strategy in (Zhao et al., 2016), in *Simulation 1* where all factor loadings are sparse, we calculate the SSI between the true and recovered factor loading matrices. In *Simulation 2*, we first threshold the recovered factor loading matrix entries with a sparsity threshold set to 0.15. Then, we categorize the columns of each recovered factor loading matrix into sparse columns and dense columns by selecting the first 4 columns with most nonzero entries as dense columns, and the remaining columns as sparse columns. We calculate SSI between the true and the recovered sparse factor loading columns, and DSI between the true and the recovered dense columns. We calculate the two stability indices for each group separately and average the result for all groups.

The true and the inferred factor loading matrices by all methods in *Simulation 1* and *Simulation 2* are shown in Fig. 7. The ARD prior cannot induce sufficient sparsity by pushing irrelevant factor loadings to small values. As a consequence, the GFA has difficulty in recovering sparse factor loadings because of the columns-wise ARD priors (Fig. 7). Similarly, the sGFA cannot induce sufficient element-wise sparsity within the loading columns by the independent ARD priors (Fig. 7). The ssGFA overfitted to data by not sufficiently shutting off the redundant factors

(Fig. 7). Both the BASS and NGFA achieve element-wise sparsity effectively (Fig. 7). We quantify the performance of the methods with stability indices, i.e., the means and the standard derivations of the stability indices for each method over 20 runs are shown in Fig. 4. The NGFA using our CVI algorithm achieves the best SSI and DSI scores almost for all sample sizes.

## 4.2 Cancer gene prioritization

Integrative analysis of multiple genomic data sets for understanding the genetic basis of common diseases has been challenging. For instance, DNA alterations that are frequent in cancers, measured by copy number variation (CNV) data, are known to induce gene expression modifications. Hence, cancer-related genes can be discovered by searching for such interactions. Recently, Bayesian GFA methods were applied to the task of cancer gene prioritization with encouraging results (Klami et al., 2013). To demonstrate the effectiveness of the NGFA using our CVI algorithm, we choose the same datasets `Hyman` and `Pollack` from (Lahti et al., 2013) that are based on gene expression (GE) and CNV data as described in Table 2.

| Dataset | # genes | # samples | # cancer genes |
|---------|---------|-----------|----------------|
| Hyman   | 7489    | 14        | 48             |
| Pollack | 4287    | 41        | 38             |

Table 2: The details of cancer genomics datasets.

More specifically, we consider the patients as co-occurring samples and all the genes in the whole genome as features. The GE and CNV data constitute the two groups. We then rank the genes according to the quantity defined by $s_d = \sum_{k=1}^{K} |\mathsf{E}(g_{kd}^{(1)}) \mathsf{E}(g_{kd}^{(2)})|$, that is the correlation between GE and CNV data captured by the shared factors. We repeat the data pre-processing procedure in

(Lahti et al., 2013), and evaluate the model performance by the area under the curve of the receiver operating characteristic (AUC) for retrieving known cancer-related genes. We run the NGFA 20 times with the initial $K$ set to the minimum of the sample size and feature dimension. We compare the NGFA using CVI algorithm to the Bayesian inter-battery factor analysis (BIBFA) model. We run the BIBFA for 20 times according to the setting described in (Klami et al., 2013). The mean AUC scores and the standard deviations are shown in Fig. 5. The AUC scores for all the other methods are cited from (Lahti et al., 2013) where the standard deviations cannot be presented because those alternatives are deterministic methods. The NGFA using our CVI algorithm outperforms all the alternative methods.

### 4.3 Decoding fMRI brain activity

Bayesian canonical correlation analysis (BCCA) was investigated to analyze fMRI responses to visual stimuli in (Fujiwara et al., 2009). We evaluate the NGFA using our CVI algorithm to the fMRI recordings of two subjects viewing visual images consisting of contrast-defined $10 \times 10$ patches (Miyawaki et al., 2008). The data is composed of two independent sessions: one for "random image session" with spatially random patterns sequentially presented; the other for "figure image session" with alphabet letters and geometric shapes sequentially presented. For the NGFA, we first treat the random image session and corresponding fMRI recordings as two groups, to extract the image bases and weight vector automatically from the input, with the initial $K$ set to $\min(D_1, D_2) = 100$. Our task is to reconstruct the visual image from the new fMRI recordings in the figure image session. The reconstruction performance is evaluated by the mean squared error between the presented and reconstructed images. We run both the BCCA and the NGFA for 20 times. The mean squared prediction error over 20 runs for NGFA is 0.224 with the standard deviation less than 1e-3, which is better than the result 0.251(0.002) of the BCCA. The reconstructed geometric shapes and alphabet letters by the BCCA and the proposed NGFA are shown in Fig. 6.
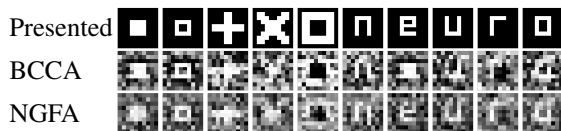


Figure 6: Presented images (first row) and the reconstructed visual images obtained from the BCCA (second row) and the NGFA (third row).

## 5 Discussion

In this work, the GFA problem is tackled via a Bayesian nonparametric method that allows the total number of factors to be automatically inferred, and the underlying structured sparsity to be effectively captured. In particular, we have presented an efficient collapsed variational inference algorithm for the nonparametric Bayesian group factor analysis model. By integrating out the group-specific beta process parameters, our CVI algorithm achieves a better approximation because all latent variables are dependent through the field while the weak dependences are very small in the collapsed space. Using the Gaussian approximation technique, all the variational parameters can be efficiently updated through closed form expressions. Experimental results on both synthetic data and real-world applications demonstrate superior performance of our CVI algorithm for the nonparametric Bayesian group factor analysis model when compared to state-of-the-art GFA methods. An interesting direction of future research is how to infer hierarchically structured latent factors, as was done for deep factor modelling (Gan et al., 2015a,b; Zhou et al., 2016). Another possible direction would be to generalize GFA methods to model dynamic multiple related graph data (Durante et al., 2017) under the Poisson factorization framework (Zhou, 2015; Yang and Koeppl, 2018).

## References

Blei, D. M. et al. (2003). Latent Dirichlet allocation. *JMLR*, 3:993–1022.

Bunte, K. et al. (2016). Sparse group factor analysis for biclustering of multiple data sources. *Bioinformatics*.

Carvalho, C. M. et al. (2008). High-dimensional sparse factor modeling: Applications in gene expression genomics. *JASA*, 103(484):1438–1456.

Chen, B. et al. (2011). The hierarchical beta process for convolutional factor analysis and deep learning. In *ICML*, pages 361–368.

Durante, D. et al. (2017). Bayesian learning of dynamic multilayer networks. *JMLR*, 18(1):1414–1442.

Foulds, J. et al. (2013). Stochastic collapsed variational
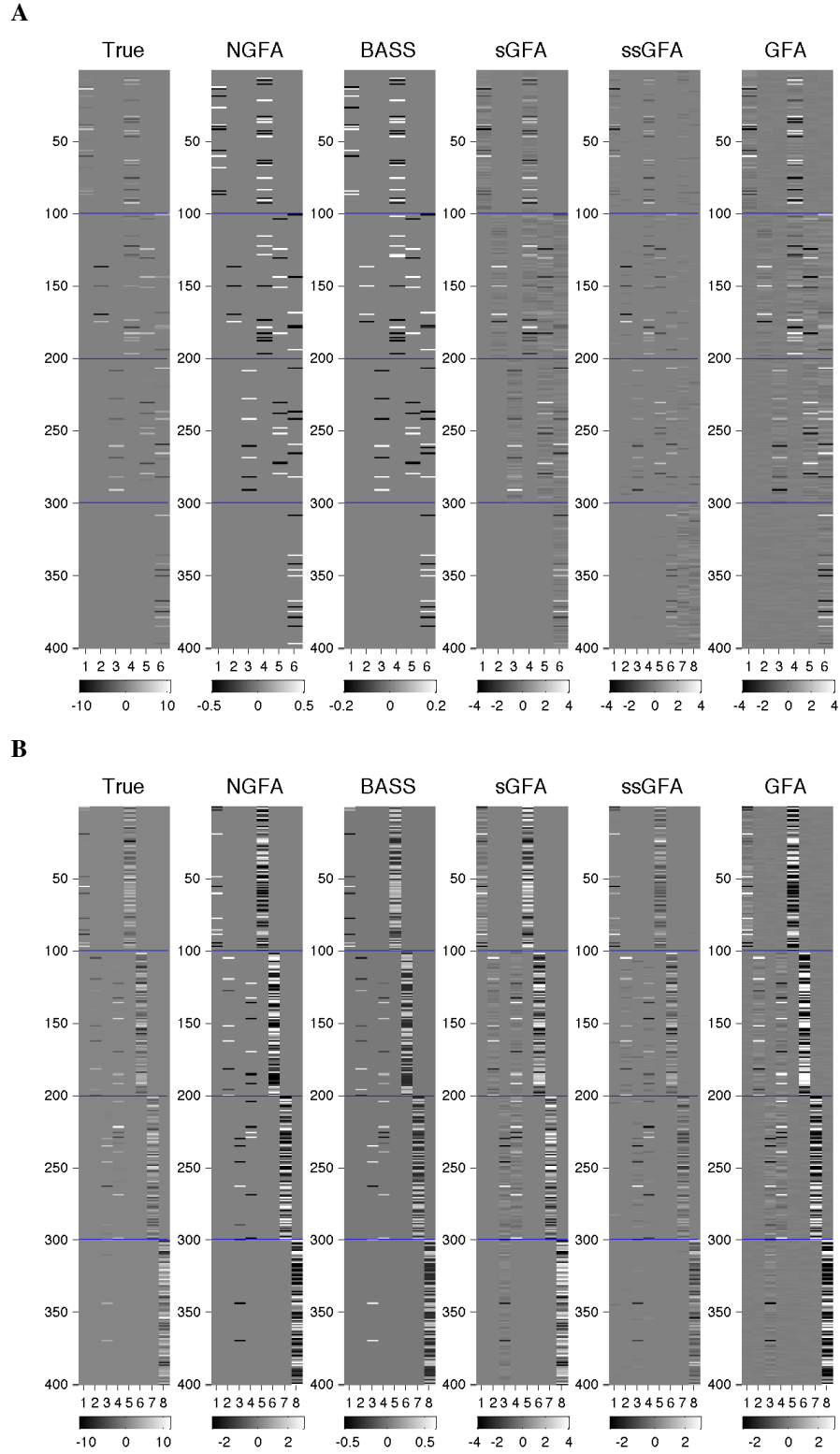
Figure 7: The true and the inferred factor loadings by all methods in *Simulation 1* (A) and *Simulation 2* (B). The columns of the inferred factor loading matrices were reordered for easy comparison. The horizontal lines separate the four groups.

Bayesian inference for latent Dirichlet allocation. In *KDD*, pages 446–454.

Fox, E. B. et al. (2011). A sticky hdp-hmm with application to speaker diarization. *Ann. Appl. Stat.*, 5:1020–1056.

Fujiwara, Y. et al. (2009). Estimating image bases for visual image reconstruction from human rain activity. In *NIPS*, pages 576–584.

Gan, Z. et al. (2015a). Learning Deep Sigmoid Belief Networks with Data Augmentation. In *AISTATS*, pages 268–276.

Gan, Z. et al. (2015b). Scalable deep Poisson factor analysis for topic modeling. In *ICML*, pages 1823–1832.

Gupta, S. K. et al. (2012a). A Bayesian nonparametric joint factor model for learning shared and individual subspaces from multiple data sources. In *SDM*, pages 200–211.

Gupta, S. K. et al. (2012b). A slice sampler for restricted hierarchical beta process with applications to shared subspace learning. In *UAI*, pages 316–325.

Hoef, J. M. V. (2012). Who invented the delta method? *The American Statistician*, 66(2):124–127.

Ishiguro, K. et al. (2017). Averaged collapsed variational Bayes inference. *JMLR*, 18(1):1–29.

Klami, A., Virtanen, S., and Kaski, S. (2013). Bayesian canonical correlation analysis. *JMLR*, 14(1):965–1003.

Knowles, D. A. et al. (2011). Nonparametric Bayesian sparse factor models with application to gene expression modeling. *Ann. Appl. Stat.*, 5(2B):1534–1552.

Lahti, L. et al. (2013). Cancer gene prioritization by integrative analysis of mRNA expression and DNA copy number data. *Briefings in Bioinformatics*, 14(1):27–35.

Miyawaki, Y. et al. (2008). Visual Image Reconstruction from Human Brain Activity using a Combination of Multiscale Local Image Decoders. *Neuron*, pages 5–29.

Paisley, J. et al. (2009). Nonparametric factor analysis with beta process priors. In *ICML*, pages 777–784.

Pitman, J. (2006). *Combinatorial stochastic processes*. Springer-Verlag, Berlin. Lectures on Probability Theory.

Rai, P. and Daume III, H. (2008). The infinite hierarchical factor regression model. In *NIPS*, pages 1321–1328.

Teh, Y. W. et al. (2006). A collapsed variational Bayesian inference algorithm for latent Dirichlet allocation. In *NIPS*, pages 1353–1360.

Teh, Y. W. et al. (2007). Hierarchical Dirichlet processes. *JASA*, 101(476):1566–1581.

Teh, Y. W. et al. (2008). Collapsed variational inference for HDP. In *NIPS*, pages 1481–1488.

Thibaux, R. et al. (2007). Hierarchical beta processes and the Indian buffet process. In *AISTATS*, pages 564–571.

Virtanen, S. et al. (2012). Bayesian group factor analysis. In *AISTATS*, pages 1269–1277.

Wainwright, M. J. et al. (2008). Graphical models, exponential families, and variational inference. *Found. Trends Mach. Learn.*

Wang, P. et al. (2013). Collapsed variational Bayesian inference for hidden markov models. In *AISTATS*, pages 599–607.

West, M. (2003). Bayesian factor regression models in the "large p, small n" paradigm. In *Bayesian Statistics*, pages 723–732.

Witten, D. M. et al. (2009). A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*, 10(3):515–534.

Yang, S. and Koeppl, H. (2018). Dependent relational gamma process models for longitudinal networks. In *ICML*, pages 5551–5560.

Zhang, B. et al. (2017). Collapsed variational bayes for Markov jump processes. In *NIPS*, pages 3749–3757.

Zhao, S. et al. (2016). Bayesian group factor analysis with structured sparsity. *JMLR*, 17(196):1–47.

Zhou, M. (2015). Infinite edge partition models for overlapping community detection and link prediction. In *AISTATS*, pages 1135–1143.

Zhou, M. et al. (2015). Negative binomial process count and mixture modeling. *IEEE Trans. PAMI*, 37(2):307–320.

Zhou, M. et al. (2016). Augmentable gamma belief networks. *JMLR*, 17:1–44.

Zou, H. et al. (2006). Sparse principal component analysis. *J. Comput. and Graph. Statist.*, 15(2):265–286.

# Appendix: Collapsed Variational Inference for Nonparametric Bayesian Group Factor Analysis

## A.2 Variational Approximation

Here, we provide the full variational approximation used in our CVI algorithm for the NGFA. We use "·" as a index summation shorthand, e.g., $x_{\cdot j} = \sum_i x_{ij}$. We assume the variational posterior over the latent variables and parameters as

$$q(\mathbf{Z}, \mathbf{W}, \mathbf{F}, \boldsymbol{\beta}, \boldsymbol{\lambda}, \boldsymbol{\tau}, \boldsymbol{\alpha}, \mathbf{s}, \mathbf{t}, \boldsymbol{\eta})$$
$$= q(\mathbf{W})q(\mathbf{F})q(\boldsymbol{\beta})q(\boldsymbol{\lambda})q(\boldsymbol{\tau})q(\boldsymbol{\alpha})q(\mathbf{s}, \mathbf{t}, \boldsymbol{\eta}|\mathbf{Z})q(\mathbf{Z}),$$

where we define the variational posterior for each parameter as

$$q(\mathbf{W}) = \prod_{m,d,k} \mathcal{N}(w_{kd}^{(m)}; \mu_{w_{kd}}^{(m)}, \sigma_{w_{kd}}^{(m)}),$$

$$q(\mathbf{F}) = \prod_{n,k} \mathcal{N}(f_{nk}; \mu_{f_{nk}}, \sigma_{f_{nk}}),$$

$$q(\boldsymbol{\beta}) = \prod_{k} \mathrm{Beta}(\beta_k; a_k, b_k),$$

$$q(\boldsymbol{\lambda}) = \prod_{m,d,k} \mathrm{Gam}(\lambda_{kd}^{(m)}; e_{kd}^{(m)}, f_{kd}^{(m)}),$$

$$q(\boldsymbol{\tau}) = \prod_{m,n} \mathrm{Gam}(\tau_n^{(m)}; g_n^{(m)}, h_n^{(m)}),$$

$$q(\boldsymbol{\alpha}) = \prod_{m} \mathrm{Gam}(\alpha^{(m)}; c^{(m)}, d^{(m)}),$$

$$q(\mathbf{Z}) = \prod_{m,d,k} \mathrm{Bern}(z_{kd}^{(m)}; \rho_{kd}^{(m)}),$$

$$q(\mathbf{s}|\mathbf{Z}) = \prod_{m,k} \begin{bmatrix} \hat{n}_{mk} \\ s_{mk} \end{bmatrix} (\mathsf{G}[\alpha^{(m)}\beta_k])^{s_{mk}},$$

$$q(\mathbf{t}|\mathbf{Z}) = \prod_{m,k} \begin{bmatrix} \tilde{n}_{mk} \\ t_{mk} \end{bmatrix} (\mathsf{G}[\alpha^{(m)}(1-\beta_k)])^{t_{mk}},$$

$$q(\boldsymbol{\eta}|\mathbf{Z}) = \prod_{m} \mathrm{Beta}(\eta_m; \mathsf{E}\left[\alpha^{(m)}\right], D_m).$$

## A.3 Evidence Lower Bound (ELBO)

The log marginal likelihood of data is lower bounded as

$$\log p(\mathbf{X} \mid \kappa_0) \geq \mathsf{E}\left[p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}, \mathbf{s}, \mathbf{t}, \boldsymbol{\eta} \mid \kappa_0)\right] - \mathsf{E}\left[q(\mathbf{Z}, \boldsymbol{\theta}, \mathbf{s}, \mathbf{t}, \boldsymbol{\eta})\right]$$
$$= \mathsf{E}_{q(\theta, \mathbf{Z})}\left[\mathsf{E}_{q(\mathbf{s}, \mathbf{t}, \boldsymbol{\eta}|\mathbf{Z})}\left[\log \frac{p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}, \mathbf{s}, \mathbf{t}, \boldsymbol{\eta} \mid \kappa_0)}{q(\mathbf{s}, \mathbf{t}, \boldsymbol{\eta} \mid \mathbf{Z})}\right] - \log q(\theta, \mathbf{Z})\right]$$
$$= \mathsf{E}_{q(\boldsymbol{\theta}, \mathbf{Z})}\left[\log p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\theta} \mid \kappa_0) - q(\boldsymbol{\theta}, \mathbf{Z})\right], \tag{12}$$

where the second equality holds provided that $q(\mathbf{s}, \mathbf{t}, \boldsymbol{\eta} \mid \mathbf{Z})$ is set to its true posterior.

To derive the variational update for each parameter, we expand the ELBO for each term in Eq. 12 as

$$
\begin{aligned}
\log p(\mathbf{X} \mid \kappa_0) \geq\ & \mathsf{E}\left[\log p(\mathbf{X}|\mathbf{W}, \mathbf{Z}, \mathbf{F}, \boldsymbol{\tau})\right] \\
& + \mathsf{E}\left[\log p(\mathbf{W})\right] - \mathsf{E}\left[\log q(\mathbf{W})\right] + \mathsf{E}\left[\log p(\mathbf{Z})\right] - \mathsf{E}\left[\log q(\mathbf{Z})\right] \\
& + \mathsf{E}\left[\log p(\mathbf{F})\right] - \mathsf{E}\left[\log q(\mathbf{F})\right] + \mathsf{E}\left[\log p(\boldsymbol{\lambda})\right] - \mathsf{E}\left[\log q(\boldsymbol{\lambda})\right] \\
& + \mathsf{E}\left[\log p(\boldsymbol{\tau})\right] - \mathsf{E}\left[\log q(\boldsymbol{\tau})\right] + \mathsf{E}\left[\log p(\boldsymbol{\alpha})\right] - \mathsf{E}\left[\log q(\boldsymbol{\alpha})\right] \\
& + \mathsf{E}\left[\log p(\boldsymbol{\beta})\right] - \mathsf{E}\left[\log q(\boldsymbol{\beta})\right].
\end{aligned}
\tag{13}
$$

### A.4 Variational Updates

The variational updates for each parameter are obtained by taking the derivate of the ELBO in Eq. 13 w.r.t. each parameter and setting it to zero.

**Updates for the sufficient statistics:**

$$
\mathsf{E}\left[\hat{n}_{mk}\right] = \sum_d \rho_{kd}^{(m)}, \qquad \mathsf{E}\left[\tilde{n}_{mk}\right] = \sum_d (1 - \rho_{kd}^{(m)}),
$$

$$
p_+ \hat{n}_{mk}) = 1 - \exp\left(\sum_d \log[1 - \rho_{kd}^{(m)}]\right),
$$

$$
p_+(\tilde{n}_{mk}) = 1 - \exp\left(\sum_d \log[\rho_{kd}^{(m)}]\right),
$$

$$
\mathsf{E}_+[\hat{n}_{mk}] = \frac{\mathsf{E}[\hat{n}_{mk}]}{p_+(\hat{n}_{mk})}, \qquad \mathsf{E}_+[\tilde{n}_{mk}] = \frac{\mathsf{E}[\tilde{n}_{mk}]}{p_+(\tilde{n}_{mk})},
$$

$$
\mathsf{V}\left[\hat{n}_{mk}\right] = \mathsf{V}\left[\tilde{n}_{mk}\right] = \sum_d (1 - \rho_{kd}^{(m)})\rho_{kd}^{(m)},
$$

$$
\mathsf{V}_+[\hat{n}_{mk}] = \frac{\mathsf{V}[\hat{n}_{mk}]}{p_+(\hat{n}_{mk})}, \qquad \mathsf{V}_+[\tilde{n}_{mk}] = \frac{\mathsf{V}[\tilde{n}_{mk}]}{p_+(\tilde{n}_{mk})}.
\tag{14}
$$

**Updates for $\sigma_{w_{kd}}^{(m)}$ and $\mu_{w_{kd}}^{(m)}$:**

$$
\sigma_{w_{kd}}^{(m)} = \left(\mathsf{E}\left[\lambda_{kd}^{(m)}\right] + \mathsf{E}\left[z_{kd}^{(m)}\right] \sum_n \mathsf{E}\left[\tau_n^{(m)}\right] \mathsf{E}\left[f_{nk}^2\right]\right)^{-1},
\tag{15}
$$

$$
\mu_{w_{kd}}^{(m)} = \sigma_{w_{kd}}^{(m)} \left(\mathsf{E}\left[z_{kd}^{(m)}\right] \sum_n \mathsf{E}\left[\tau_n^{(m)}\right] \mathsf{E}\left[f_{nk}\right] \tilde{x}_{nd}^{(m)\,-k}\right).
\tag{16}
$$

**Updates for the auxiliary variables $\mathsf{s}, \mathsf{t}$:**

$$
\begin{aligned}
\mathsf{E}[s_{mk}] \approx\ & \mathsf{G}[\alpha^{(m)}\beta_k]p_+(\hat{n}_{mk})\big(\Psi\big(\mathsf{G}[\alpha^{(m)}\beta_k] + \mathsf{E}_+[\hat{n}_{mk}]\big) \\
& - \Psi(\mathsf{G}(\alpha^{(m)}\beta_k)) + \frac{\mathsf{V}_+[\hat{n}_{mk}]\Psi'(\mathsf{G}[\alpha^{(m)}\beta_k] + \mathsf{E}_+[\hat{n}_{mk}])}{2}\big),
\end{aligned}
$$

$$
\begin{aligned}
\mathsf{E}[t_{mk}] \approx\ & \mathsf{G}[\alpha^{(m)}\bar{\beta}_k]p_+(\tilde{n}_{mk})\big(\Psi\big(\mathsf{G}[\alpha^{(m)}\bar{\beta}] + \mathsf{E}_+[\tilde{n}_{mk}]\big) \\
& - \Psi(\mathsf{G}(\alpha^{(m)}\beta_k)) + \frac{\mathsf{V}_+[\tilde{n}_{mk}]\Psi'(\mathsf{G}[\alpha^{(m)}\bar{\beta}] + \mathsf{E}_+[\tilde{n}_{mk}])}{2}\big).
\end{aligned}
\tag{17}
$$

**Updates for $\sigma_{f_{nk}}$ and $\mu_{f_{nk}}$:**

$$
\sigma_{f_{nk}} = \left(\sum_{m,d} \mathsf{E}\left[\tau_n^{(m)}\right] \mathsf{E}\left[z_{kd}^{(m)}\right] \mathsf{E}\left[\left(w_{kd}^{(m)}\right)^2\right] + 1\right)^{-1},
\tag{18}
$$

$$
\mu_{f_{nk}} = \sigma_{f_{nk}} \left(\sum_{m,d} \mathsf{E}\left[\tau_n^{(m)}\right] \mathsf{E}\left[z_{kd}^{(m)}\right] \mathsf{E}\left[w_{kd}^{(m)}\right] \tilde{x}_{nd}^{(m)\,-k}\right).
\tag{19}
$$

**Updates for $a_k$ and $b_k$:**

$$
a_k = \kappa_0/K + \mathsf{E}\left[s_{\cdot k}\right], \ b_k = \kappa_0(1 - 1/K) + \mathsf{E}\left[t_{\cdot k}\right].
\tag{20}
$$

**Algorithm 1:** Collapsed variational inference for the NGFA

---
**Input** : Data $\mathbf{X}$, Model $\log p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}, \mathbf{s}, \mathbf{t}, \boldsymbol{\eta})$, maximum iteration $\mathcal{J}$, variational approximation $q(\mathbf{Z}, \boldsymbol{\theta}, \mathbf{s}, \mathbf{t}, \boldsymbol{\eta}; \boldsymbol{\Phi})$, and hyper-parameter $\kappa_0$

**Output** : Variational parameters $\boldsymbol{\Phi}^3$

Initialize $\boldsymbol{\Phi}$ randomly.

**for** $iter = 1 : \mathcal{J}$ **do**

    **for** $k = 1$ *to* $K_+{}^4$ **do**

        Update $a_k, b_k$ (Eq. 20)

        **for** $m = 1$ *to* $M$ **do**

            Update the sufficient statistics in (Eq. 14)

            Calculate $\mathsf{E}[s_{mk}], \mathsf{E}[t_{mk}]$ (Eq. 17)

            **for** $d = 1$ *to* $D_m$ **do**

                Update $\rho_{kd}^{(m)}$ (Eq. 8) in Fig. 3

                Update $\sigma_{w_{kd}}^{(m)}, \mu_{w_{kd}}^{(m)}$ (Eq. 15; 16)

                Update $e_{kd}^{(m)}$ and $f_{kd}^{(m)}$ (Eq. 21)

            **end**

        **end**

        **for** $n = 1$ *to* $N$ **do**

            Update $\sigma_{f_{kn}}$ and $\mu_{f_{kn}}$ (Eq. 18; 19)

        **end**

    **end**

    **for** $m = 1$ *to* $M$ **do**

        Update $c^{(m)}$ and $d^{(m)}$ (Eq. 23)

        Calculate $\mathsf{E}[\log \eta_m]$ (Eq. 24)

        **for** $n = 1$ *to* $N$ **do**

            Update $g_n^{(m)}$ and $h_n^{(m)}$ (Eq. 22)

        **end**

    **end**

**end**

---

**Updates for $e_{kd}^{(m)}$ and $f_{kd}^{(m)}$:**

$$e_{kd}^{(m)} = e_0 + 1/2, \qquad f_{kd}^{(m)} = f_0 + \left( \mathsf{E}\left[ \left( w_{kd}^{(m)} \right)^2 \right] \right)/2. \tag{21}$$

**Updates for $g_n^{(m)}$ and $h_n^{(m)}$:**

$$g_n^{(m)} = g_0 + (D_m)/2, h_n^{(m)} = h_0 + \left( \mathsf{E}\left[ \| \mathbf{x}_n^{(m)} - \mathbf{G}^{(m)} \mathbf{f}_n \|^2 \right] \right)/2. \tag{22}$$

**Updates for $c^{(m)}$ and $d^{(m)}$:**

$$c^{(m)} = c_0 + \mathsf{E}\left[ s_{m\cdot} \right] + \mathsf{E}\left[ t_{m\cdot} \right], d^{(m)} = d_0 - \mathsf{E}\left[ \log \eta_m \right]. \tag{23}$$

**Updates for the auxiliary variables $\eta$:**

$$\mathsf{E}[\log \eta_m] = \Psi(\mathsf{E}[\alpha^{(m)}]) - \Psi(\mathsf{E}[\alpha^{(m)}] + D_m). \tag{24}$$

Altogether, our CVI algorithm for the NGFA is summarized in Algorithm 1.

---

[4]For the sake of clarity, we use $\boldsymbol{\Phi}$ to denote all the variational parameters.

[4]We use $K_+$ to denote the number of active factors as the hierarchical beta Bernoulli prior can shrink the coefficients of the redundant factors to zeros.