

## 0.1 Hard sparsity prior: spike and Slab

### 0.1.1 The model

We have a data set  $\mathbf{Y}$  of  $M$  input views with dimensionality  $\mathbf{Y}^m \in \mathbb{R}^{N \times D_m}$ . The sparse group factor analysis model is defined as follows:

$$y_{nd} \approx \mathcal{N}(y_{nd}^m | \mathbf{w}_{d,:}^m \mathbf{z}_{n,:}, 1/\tau_d^m) \quad (1)$$

$$w_{dk}^m \approx \theta \mathcal{N}(w_{dk}^m | 0, 1/\alpha_k^m) + (1 - \theta) \delta_0(w_{dk}^m) \quad (2)$$

$$z_{nk} \approx \mathcal{N}(z_{nk} | 0, 1) \quad (3)$$

### 0.1.2 Efficient variational inference

The presence of the Dirac delta mass function makes the application of variational inference troublesome. However, there exists a simple reparameterization of the spike and slab prior that is more amenable to approximate inference.

Assume a Gaussian random variable  $w \approx \mathcal{N}(\hat{w} | 0, \sigma^2)$  and a Bernoulli random variable  $s \approx \pi^s (1 - \pi)^{1-s}$ . The product  $s\hat{w}$  forms a new random variable distributed according to the original spike and slab prior. Thus, we can reparametrise  $w = s\hat{w}$  and assign the prior distributions on  $s$  and  $\hat{w}$ . Thus, the reparametrised spike and slab takes the form:

$$p(\hat{w}, s) = \mathcal{N}(\hat{w} | 0, \sigma^2) \pi^s (1 - \pi)^{1-s}$$

A simple approach for variational inference is to consider the mean-field approximation:

$$q(\hat{w}, s) = q(\hat{w})q(s)$$

which leads to easy updates. However, each pair of variables  $\{s\hat{w}\}$  is strongly correlated since their product is the underlying variable that interacts with the data. Therefore, considering a unimodal distribution with the variables  $\hat{w}$  and  $s$  being independent leads to a very inefficient inference. Instead, we introduce a paired mean field approximation  $q(w, s)$  which approximates better the factorial nature of the true posterior.

$$q(\hat{\mathbf{W}}, \mathbf{S}) = \prod_{m=1}^M \prod_{d=1}^{D_m} \prod_{k=1}^K q(\hat{w}_{dk}^m, s_{dk}^m)$$

Joint probability density function:

$$\begin{aligned} p(\mathbf{Y}, \hat{\mathbf{W}}, \mathbf{S}, \mathbf{Z}) &= \prod_{m=1}^M \prod_{n=1}^N \prod_{d=1}^{D_m} \mathcal{N}\left(y_{nd}^m | \sum_{k=1}^K s_{dk}^m \hat{w}_{dk}^m z_{nk}, 1/\tau_d\right) \times \\ &\quad \prod_{d=1}^{D_m} \prod_{k=1}^K \mathcal{N}(\hat{w}_{dk}^m | 0, 1/\alpha_k^m) \theta^{s_{dk}^m} (1 - \theta)^{1-s_{dk}^m} \times \\ &\quad \prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(z_{nk} | 0, 1) \end{aligned}$$

Full variational distribution:

$$q(\hat{\mathbf{W}}, \mathbf{S}, \mathbf{Z}) = \prod_{m=1}^M \prod_{d=1}^{D_m} \prod_{k=1}^K q(\hat{w}_{dk}^m, s_{dk}^m) \prod_{k=1}^K \prod_{n=1}^N q(z_{n,k})$$

### 0.1.3 Derivation of update equations

The optimal distribution  $\hat{q}_i$  for each variable  $\mathbf{x}_i$ , is the following:

$$\log \hat{q}_i(\mathbf{x}_i) = \mathbb{E}_{i \neq j}[\log p(\mathbf{Y}, \mathbf{X})] + \text{const.} \quad (4)$$

where  $\mathbb{E}_{i \neq j}$  denotes an expectation with respect to the  $q$  distributions over all variables  $\mathbf{x}_j$  except for  $\mathbf{x}_i$ . The additive constant is set by normalising the distribution  $\hat{q}_i(\mathbf{z}_i)$ :

$$\hat{q}_i(\mathbf{x}_i) = \frac{\exp(\mathbb{E}_{i \neq j}[\log p(\mathbf{Y}, \mathbf{X})])}{\int \exp(\mathbb{E}_{i \neq j}[\log p(\mathbf{Y}, \mathbf{X})]) d\mathbf{X}}$$

### Latent variables

Term from the likelihood  $p(\mathbf{Y}|\hat{\mathbf{W}}, \mathbf{Z}, \boldsymbol{\tau}, \mathbf{S})$ :

$$\begin{aligned} & \sum_{m=1}^M \sum_{d=1}^{D_m} \langle \tau_d^m \rangle \langle s_{dk}^m \hat{w}_{dk}^m \rangle y_{nd}^m z_{nk} - \frac{1}{2} \sum_{m=1}^M \sum_{d=1}^{D_m} \langle \tau_d^m \rangle \langle (s_{dk}^m \hat{w}_{dk}^m)^2 \rangle z_{nk}^2 \\ & - \frac{1}{2} \sum_{m=1}^M \frac{1}{2} \sum_{d=1}^D \langle \tau_d^m \rangle \sum_{j \neq k} (\langle s_{dj}^m \hat{w}_{dj}^m \rangle \langle z_{nj} \rangle) \langle \hat{w}_{dk}^m s_{dk}^m \rangle \langle z_{nk} \rangle + \text{const.} \end{aligned}$$

Term from the prior  $p(z_{nk})$ :

$$-\frac{1}{2} z_{nk}^2 + \text{const.}$$

Variational distribution:

$$q(\mathbf{Z}) = \prod_{k=1}^K \prod_{n=1}^N q(z_{nk}) = \prod_{k=1}^K \prod_{n=1}^N \mathcal{N}(z_{nk} | \mu_{z_{nk}}, \sigma_{z_{nk}})$$

where

$$\begin{aligned} \sigma_{z_{nk}}^2 &= \left( \sum_{m=1}^M \sum_{d=1}^{D_m} \tau_d^m \langle (s_{dk}^m \hat{w}_{dk}^m)^2 \rangle + 1 \right)^{-1} \\ \mu_{z_{nk}} &= \sigma_{z_{nk}}^2 \sum_{m=1}^M \sum_{d=1}^{D_m} \langle \tau_d^m \rangle \langle s_{dk}^m \hat{w}_{dk}^m \rangle \left( y_{nd}^m - \sum_{j \neq k} \langle s_{dj}^m \hat{w}_{dj}^m \rangle \langle z_{nj} \rangle \right) \end{aligned}$$

### Spike and Slab Weights

Variational distribution:

$$q(\hat{\mathbf{W}}, \mathbf{S}) = \prod_{m=1}^M \prod_{d=1}^{D_m} \prod_{k=1}^K q(\hat{w}_{dk}^m, s_{dk}^m) = \prod_{m=1}^M \prod_{d=1}^{D_m} \prod_{k=1}^K q(\hat{w}_{dk}^m | s_{dk}^m) q(s_{dk}^m)$$

Update for  $q(s_{dk}^m)$ :

$$\gamma_{dk} = q(s_{dk} = 1) = \frac{1}{1 + \exp(-\lambda_{dk})}$$

where

$$\lambda_{dk}^m = \langle \log \frac{\theta}{1-\theta} \rangle + 0.5 \log \frac{\langle \alpha_k^m \rangle}{\langle \tau_d^m \rangle} - 0.5 \log \left( \sum_{n=1}^N \langle z_{nk}^2 \rangle + \frac{\langle \alpha_k^m \rangle}{\langle \tau_d^m \rangle} \right) + \frac{\langle \tau_d^m \rangle}{2} \frac{\left( \sum_{n=1}^N y_{nd}^m \langle z_{nk} \rangle - \sum_{j \neq k} \langle s_{dj}^m \hat{w}_{dj}^m \rangle \sum_{n=1}^N \langle z_{nk} \rangle \langle z_{nj} \rangle \right)^2}{\sum_{n=1}^N \langle z_{nk}^2 \rangle + \frac{\langle \alpha_k^m \rangle}{\langle \tau_d^m \rangle}}$$

Update for  $q(\hat{w}_{dk}^m)$ :

$$q(\hat{w}_{dk}^m | s_{dk}^m = 0) = \mathcal{N}(\hat{w}_{dk}^m | 0, 1/\alpha_k^m)$$

$$q(\hat{w}_{dk}^m | s_{dk}^m = 1) = \mathcal{N}(\hat{w}_{dk}^m | \mu_{w_{dk}}^m, \sigma_{w_{dk}}^2)$$

where

$$\mu_{w_{dk}}^m = \frac{\sum_{n=1}^N y_{nd}^m \langle z_{nk} \rangle - \sum_{j \neq k} \langle s_{dj}^m \hat{w}_{dj}^m \rangle \sum_{n=1}^N \langle z_{nk} \rangle \langle z_{nj} \rangle}{\sum_{n=1}^N \langle z_{nk}^2 \rangle + \frac{\langle \alpha_k^m \rangle}{\langle \tau_d^m \rangle}}$$

$$\sigma_{w_{dk}}^m = \frac{\langle \tau_d^m \rangle^{-1}}{\sum_{n=1}^N \langle z_{nk}^2 \rangle + \frac{\langle \alpha_k^m \rangle}{\langle \tau_d^m \rangle}}$$

Taken together this means that we can update  $q(\hat{w}_{dk}^m, s_{dk}^m)$  using:

$$q(\hat{w}_{dk}^m | s_{dk}^m) \times q(s_{dk}^m) = \mathcal{N}(\hat{w}_{dk}^m | s_{dk}^m \mu_{w_{dk}}^m, s_{dk}^m \sigma_{w_{dk}}^2 + (1 - s_{dk}^m)/\alpha_k^m) \times (\lambda_{dk}^m)^{s_{dk}^m} (1 - \lambda_{dk}^m)^{1-s_{dk}^m}$$

## ARD precision (alpha)

$$\log \hat{q}_i(\mathbf{x}_i) = \mathbb{E}_{i \neq j} [\log p(\mathbf{Y}, \mathbf{X})] + \text{const.} \quad (5)$$

Term from the prior  $\log p(\alpha_k^m)$ :

$$(a_0^\alpha - 1) \log(\alpha_k^m) - b_0^\alpha \alpha_k^m + \text{const.}$$

Term from the prior  $\log p(\mathbf{w}_{:k}^m) = \sum_{d=1}^{D_m} \log p(s_{dk}^m, \hat{w}_{dk}^m)$ :

$$\frac{D_m}{2} \log(\alpha_k^m) - \frac{\alpha_k^m}{2} \sum_{d=1}^D \langle \hat{w}_{dk}^2 \rangle + \sum_{d=1}^{D_m} \{ \langle s_{dk}^m \rangle \log \theta_0 + (1 - \langle s_{dk}^m \rangle) \log(1 - \theta_0) \} + \text{const.}$$

Writing everything together:

$$\left( a_0^\alpha + \frac{D_m}{2} - 1 \right) \log \alpha_k^m - \left( b_0^\alpha + \frac{1}{2} \sum_{d=1}^{D_m} \langle (\hat{w}_{d,k}^m)^2 \rangle \right) \alpha_k^m + \text{const.}$$

Variational distribution:

$$q(\boldsymbol{\alpha}) = \prod_{m=1}^M \prod_{k=1}^K \mathcal{G}(\alpha_k^m | \hat{a}_{mk}^\alpha, \hat{b}_{mk}^\alpha)$$

where

$$\hat{a}_{mk}^\alpha = a_0^\alpha + \frac{D_m}{2}$$

$$\hat{b}_{mk}^\alpha = b_0^\alpha + \frac{\sum_{d=1}^{D_m} \langle (\hat{w}_{d,k}^m)^2 \rangle}{2}$$

## Noise precision (tau)

Term from the prior  $p(\tau_d^m)$ :

$$(a_0^\tau - 1) \log \tau_d^m - b_0^\tau \tau_d^m + \text{const.}$$

Term from the likelihood  $\mathcal{N}(y_{n,d}^m | \sum_{k=1}^K \hat{w}_{dk}^m s_{dk}^m z_{nk}, \tau_d^m)$ :

$$\frac{N}{2} \log \tau_d^m - \frac{\tau_d^m}{2} \sum_{n=1}^N \langle (y_{nd}^m - \sum_k w_{dk}^m s_{dk}^m z_{nk})^2 \rangle + \text{const.}$$

Rewriting everything together:

$$\left(a_0^\tau - 1 + \frac{N}{2}\right) \log \tau_d^m - \left(b_0 + \frac{1}{2} \sum_{n=1}^N \langle (y_{nd}^m - \sum_k \hat{w}_{dk}^m s_{dk}^m z_{nk})^2 \rangle\right) \tau_d^m$$

Variational distribution:

$$q(\boldsymbol{\tau}) = \prod_{m=1}^M \prod_{d=1}^{D_m} q(\tau_d^m) = \prod_{m=1}^M \prod_{d=1}^{D_m} \mathcal{G}(\tau_d^m | \hat{a}_{md}^\tau, \hat{b}_{md}^\tau)$$

where

$$\begin{aligned} \hat{a}_{md}^\tau &= a_0^\tau + \frac{N}{2} \\ \hat{b}_{md}^\tau &= b_0^\tau + \frac{1}{2} \sum_{n=1}^N \langle (y_{nd}^m - \sum_k \hat{w}_{dk}^m s_{dk}^m z_{n,k})^2 \rangle \end{aligned}$$

## Spike and Slab sparsity parameter $\theta$

Unless a given factor is specifically annotated in a given view, the sparsity parameter  $\theta_k^m$  of the Spike and Slab prior on  $w_{k,d}^m, \forall d$  is given a Beta prior:  $P(\theta_k^m) = \text{Beta}(a_0, b_0)$ . The posterior  $q(\theta_k^m)$  is Beta distributed and the update of its parameters  $a_k^m$  and  $b_k^m$  are given below:

$$\begin{aligned} a_k^m &= \sum_d \langle S_{k,d}^m \rangle + a_0 \\ b_k^m &= b_0 - \sum_d \langle S_{k,d}^m \rangle + D_m \end{aligned}$$

### 0.1.4 Lower bound

#### Likelihood term

Vector form:

$$- \sum_{m=1}^M \frac{ND_m}{2} \log(2\pi) + \frac{N}{2} \sum_{m=1}^M \sum_{d=1}^{D_m} \log(\tau_d^m) - \frac{1}{2} \sum_{m=1}^M \sum_{d=1}^{D_m} \left( \mathbf{y}_d^m - \sum_{k=1}^K \langle s_{dk}^m \hat{w}_{dk}^m \rangle \langle \mathbf{z}_k \rangle \right)^T (\tau_d^m \mathbf{I}) \left( \mathbf{y}_d^m - \sum_{k=1}^K \langle s_{dk}^m \hat{w}_{dk}^m \rangle \langle \mathbf{z}_k \rangle \right)$$

Scalar form:

$$- \sum_{m=1}^M \frac{ND_m}{2} \log(2\pi) + \frac{N}{2} \sum_{m=1}^M \sum_{d=1}^{D_m} \log(\langle \tau_d^m \rangle) - \sum_{m=1}^M \sum_{d=1}^{D_m} \frac{\langle \tau_d^m \rangle}{2} \sum_{n=1}^N \left( y_{nd}^m - \sum_{k=1}^K \langle s_{dk}^m \hat{w}_{dk}^m \rangle \langle z_{nk} \rangle \right)^2$$

Extending terms and rearranging:

### W and S terms

$p(\hat{\mathbf{W}}, \mathbf{S})$ :

$$\begin{aligned}
& - \sum_{m=1}^M \frac{KD_m}{2} \log(2\pi) + \sum_{m=1}^M \frac{D_m}{2} \sum_{k=1}^K \log(\alpha_k^m) - \sum_{m=1}^M \frac{\alpha_k^m}{2} \sum_{d=1}^{D_m} \sum_{k=1}^K \langle (\hat{w}_{dk}^m)^2 \rangle \\
& + \langle \log(\theta) \rangle \sum_{m=1}^M \sum_{d=1}^{D_m} \sum_{k=1}^K \langle s_{dk}^m \rangle + \langle \log(1-\theta) \rangle \sum_{m=1}^M \sum_{d=1}^{D_m} \sum_{k=1}^K (1 - \langle s_{dk}^m \rangle)
\end{aligned}$$

$q(\hat{\mathbf{W}}, \mathbf{S})$ :

$$\begin{aligned}
& - \sum_{m=1}^M \frac{KD_m}{2} \log(2\pi) + \frac{1}{2} \sum_{m=1}^M \sum_{d=1}^{D_m} \sum_{k=1}^K \log(\langle s_{dk}^m \rangle \sigma_{w_{dk}^m}^2 + (1 - \langle s_{dk}^m \rangle) / \alpha_k^m) \\
& + \sum_{m=1}^M \sum_{d=1}^{D_m} \sum_{k=1}^K (1 - \langle s_{dk}^m \rangle) \log(1 - \langle s_{dk}^m \rangle) - \langle s_{dk}^m \rangle \log \langle s_{dk}^m \rangle
\end{aligned}$$

### Z term

$$\begin{aligned}
\mathbb{E}[\log P(\mathbf{Z})] &= -\frac{NK}{2} \log(2\pi) - \frac{1}{2} \sum_{n=1}^N \langle z_{nk}^2 \rangle \\
\mathbb{E}[\log q(\mathbf{Z})] &= -\frac{NK}{2} (1 + \log(2\pi)) - \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \log(\sigma_{z_{nk}}^2)
\end{aligned}$$

### alpha term

$$\begin{aligned}
\mathbb{E}[\log p(\boldsymbol{\alpha})] &= \sum_{m=1}^M \sum_{k=1}^K \left( a_0^\alpha \log b_0^\alpha + (a_0^\alpha - 1) \langle \log \alpha_k \rangle - b_0^\alpha \langle \alpha_k \rangle - \log \Gamma(a_0^\alpha) \right) \\
\mathbb{E}[\log q(\boldsymbol{\alpha})] &= \sum_{m=1}^M \sum_{k=1}^K \left( \hat{a}_k^\alpha \log \hat{b}_k^\alpha + (\hat{a}_k^\alpha - 1) \langle \log \alpha_k \rangle - \hat{b}_k^\alpha \langle \alpha_k \rangle - \log \Gamma(\hat{a}_k^\alpha) \right)
\end{aligned}$$

### tau

$$\begin{aligned}
\mathbb{E}[\log P(\boldsymbol{\tau})] &= \sum_{m=1}^M D_m a_0^\tau \log b_0^\tau + \sum_{m=1}^M \sum_{d=1}^{D_m} (a_0^\tau - 1) \langle \log \tau_d^m \rangle - \sum_{m=1}^M \sum_{d=1}^{D_m} b_0^\tau \langle \tau_d^m \rangle - \sum_{m=1}^M D_m \Gamma(a_0^\tau) \\
\mathbb{E}[\log Q(\boldsymbol{\tau})] &= \sum_{m=1}^M \sum_{d=1}^{D_m} \left( \hat{a}_{dm}^\tau \log \hat{b}_{dm}^\tau + (\hat{a}_{dm}^\tau - 1) \langle \log \tau_d^m \rangle - \hat{b}_{dm}^\tau \langle \tau_d^m \rangle - \log \Gamma(\hat{a}_{dm}^\tau) \right)
\end{aligned}$$

### Theta

$$\begin{aligned}
\mathbb{E}[\log P(\theta)] &= \sum_{m=1}^M \sum_{k=1}^K \sum_{d=1}^{D_m} ((a_0 - 1) \times \langle \log(\pi_{d,k}^m) \rangle + (b_0 - 1) \langle \log(1 - \pi_{d,k}^m) \rangle - \log(B(a_0, b_0))) \\
\mathbb{E}[\log Q(\theta)] &= \sum_{m=1}^M \sum_{k=1}^K \sum_{d=1}^{D_m} ((a_{k,d}^m - 1) \times \langle \log(\pi_{d,k}^m) \rangle + (b_{k,d}^m - 1) \langle \log(1 - \pi_{d,k}^m) \rangle - \log(B(a_{k,d}^m, b_{k,d}^m)))
\end{aligned}$$

where the expectations are calculated as follows:

## Expectations

The expectations are calculated as follows:

$$\begin{aligned}
\langle s_{dk}^m \hat{w}_{dk}^m \rangle &= \lambda_{dk}^m \mu_{w_{dk}^m} \\
\langle s_{dk}^m \hat{w}_{dk}^{m2} \rangle &= \lambda_{dk}^m (\mu_{w_{dk}^m}^2 + \sigma_{w_{dk}^m}^2) \\
\langle \hat{w}_{dk}^{m2} \rangle &= \lambda_{dk}^m (\mu_{w_{dk}^m}^2 + \sigma_{w_{dk}^m}^2) + (1 - \lambda_{dk}^m) / \alpha_k^m \\
\langle z_{nk} \rangle &= \mu_{z_{nk}} \\
\langle z_{nk}^2 \rangle &= \mu_{z_{nk}}^2 + \sigma_{z_{nk}}^2 \\
\langle \tau_d^m \rangle &= \tilde{a}_{md}^\tau / \tilde{b}_{md}^\tau \\
\langle \log \tau_d^m \rangle &= \psi(\tilde{a}_{md}^\tau) - \log \tilde{b}_{md}^\tau
\end{aligned}$$

$$\begin{aligned}
&\langle (y_{nd}^m - \sum_k^K w_{dk}^m s_{dk}^m z_{nk})^2 \rangle = \\
&(y_{nd}^m)^2 - 2y_{nd}^m \sum_k^K \langle w_{dk}^m s_{dk}^m \rangle \langle z_{nk} \rangle + \sum_k^K \sum_j^K \langle w_{dk}^m s_{dk}^m \rangle \langle z_{nk} \rangle \langle w_{dj}^m s_{dj}^m \rangle \langle z_{nj} \rangle = \\
&(y_{nd}^m)^2 - 2y_{nd}^m \sum_k^K \langle w_{dk}^m s_{dk}^m \rangle \langle z_{nk} \rangle + \sum_k^K \langle (w_{dk}^m s_{dk}^m)^2 \rangle \langle z_{nk}^2 \rangle + 2 \sum_{j>k}^K \langle w_{dk}^m s_{dk}^m \rangle \langle w_{dj}^m s_{dj}^m \rangle \langle z_{nk} \rangle \langle z_{nj} \rangle
\end{aligned}$$

## 0.2 Add cluster specific prior on the latent variables

The samples are divided into  $C$  clusters. The index  $c$  refers to one of these clusters,  $c \in [[1, C]]$ . When a sample  $n$  belongs to cluster  $c$  we write  $n \in c$  (meaning that the notation  $c$  is used both to designate the cluster and its index)

### 0.2.1 Prior architecture on $Z$

A cluster specific prior is used for every latent variable:

$$z_{n,k} \sim \mathcal{N}(\mu_{k,c}, \sigma_{k,c}), \forall n \in c, \quad (6)$$

where  $\mu_{k,c}$  is given a normal prior:

$$\mu_{k,c} \sim \mathcal{N}(\mu_0, \sigma_0), \forall k, c \quad (7)$$

### 0.2.2 New update for $Z$

$$\begin{aligned}
\sigma_{z_{nk}}^2 &= \left( \sum_{m=1}^M \sum_{d=1}^{D_m} \tau_d^m \langle (s_{dk}^m \hat{w}_{dk}^m)^2 \rangle + 1 \right)^{-1} \\
\mu_{z_{nk}} | n \in c &= \sigma_{z_{nk}}^2 \left[ \frac{\langle \mu_{k,c} \rangle}{1} + \sum_{m=1}^M \sum_{d=1}^{D_m} \langle \tau_d^m \rangle \langle s_{dk}^m \hat{w}_{dk}^m \rangle \left( y_{nd}^m - \sum_{j \neq k} \langle s_{dj}^m \hat{w}_{dj}^m \rangle \langle z_{nj} \rangle \right) \right]
\end{aligned}$$

### 0.2.3 Update for $\mu_{k,c}$

$$\sigma_{\mu_{k,c}}^2 = \left( \sum_{n \in c} \frac{1}{\sigma_{z_{n,k}}^2} + \frac{1}{\sigma_0^2} \right)^{-1}$$
$$\mu_{\mu_{k,c}} = \sigma_{\mu_{k,c}}^2 \times \left( \sum_{n \in c} \frac{\langle z_{n,k} \rangle}{\sigma_{z_{n,k}}^2} + \frac{\mu_0}{\sigma_0^2} \right)$$