

Deep Learning with Hierarchical Convolutional Factor Analysis

¹Bo Chen, ²Gungor Polatkan, ³Guillermo Sapiro, ²David Blei,
⁴David Dunson and ¹Lawrence Carin

¹Electrical & Computer Engineering Department, ⁴Statistical Sciences Department
 Duke University, Durham, NC, USA

²Computer Science Department
 Princeton University, Princeton, NJ

³Electrical & Computer Engineering Department
 University of Minnesota, Minneapolis, MN, USA

SUPPLEMENTARY MATERIAL: GIBBS UPDATE EQUATIONS

For each MCMC iteration, the samples are drawn from the following conditional distributions:

- Sampling b_{nk} and \mathbf{W}_{nk} (*i.e.*, $\{w_{nki}\}_{i \in \mathcal{S}}$): we have

$$p(b_{nk} = 1 | -) = \tilde{\pi}_{nk} \quad (1)$$

$$p(w_{nki}, i \in \mathcal{S} | -) = (1 - b_{nk})\mathcal{N}(0, \alpha_{nki}^{-1}) + b_{nk}\mathcal{N}(\mu_{nki}, \Sigma_{nki}) \quad (2)$$

where

$$\Sigma_{nki} = (\mathbf{d}_{ki}^T \mathbf{d}_{ki} \gamma_n + \alpha_{nki})^{-1} \quad (3)$$

$$\mu_{nki} = \Sigma_{nki} \gamma_n \mathbf{X}_{nki}^T \mathbf{d}_{ki} \quad (4)$$

$$\frac{\tilde{\pi}_{nk}}{1 - \tilde{\pi}_{nk}} = \frac{\pi_k}{1 - \pi_k} \cdot \frac{\mathcal{N}(\mathbf{X}_{nk} | \mathbf{W}_{nk} * \mathbf{d}_k, \gamma_n^{-1} \mathbf{I}_P)}{\mathcal{N}(\mathbf{X}_{nk} | \mathbf{0}, \gamma_n^{-1} \mathbf{I}_P)} \quad (5)$$

Here $\mathbf{X}_{nk} = \mathbf{X}_{-n} + \mathbf{W}_{nk} * \mathbf{d}_k$, $\mathbf{X}_{-n} = \mathbf{X}_n - \sum_{k=1}^K b_{nk} \mathbf{W} * \mathbf{d}_k$ and b_{nk} is the most recent sample.

Taking advantage of the convolution property, we simultaneously update the posterior mean and covariance of the coefficients for all the shifted versions of one dictionary element. Consequently,

$$\Sigma_{nk} = \mathbf{1} \oslash (\gamma_n \|\mathbf{d}_k\|_2^2 b_{nk} + \alpha_{nk}) \quad (6)$$

$$\boldsymbol{\mu}_{nk} = b_{nk} \gamma_n \Sigma_{nk} \odot (\mathbf{X}_{-n} * \mathbf{d}_k + \|\mathbf{d}_k\|_2^2 \mathbf{W}_{nk}) \quad (7)$$

where both of Σ_{nk} and $\boldsymbol{\mu}_{nk}$ have the same size with \mathbf{W}_{nk} . The symbol \odot is the element-wise product operator and \oslash the element-wise division operator.

After sampling \mathbf{W}_{nk} , we perform max-pooling on it at each Gibbs iteration (within each max-pooling region, we set all coefficients to zero, other than the maximum-amplitude coefficient). We do this for two reasons: (i) we have found that max-pooling can speed up the model inference, in that it explicitly imposes sparseness; and (ii) as stated in [?], without sparsifying the hidden units the parameter learning may be unstable (unless one employs a very large number of Gibbs samples, which is undesirable computationally). We now continue to the rest of the Gibbs draws.

- Sampling d_k :

$$P(d_{kj}|-) = \mathcal{N}(\xi_{kj}, \Lambda_{kj}) \quad (8)$$

where ξ_{kj} and Λ_{kj} is the j -th element of the following vectors ξ_k and Λ_k respectively

$$\Lambda_k = \mathbf{1} \odot \left(\sum_{n=1}^N \gamma_n b_{nk} \|\mathbf{W}_{nk}\|_2^2 + \beta_k \right), \quad (9)$$

$$\xi_k = \Lambda_k \odot \left(\sum_{n=1}^N b_{nk} \gamma_n (\mathbf{X}_{-n} * \mathbf{W}_{nk} + \mathbf{d}_k \|\mathbf{W}_{nk}\|_2^2) \right). \quad (10)$$

Here $\mathbf{X}_{nki(i)} \in \mathbb{R}^J$ represents the corresponding values of \mathbf{X}_{nki} at the same support locations of the i -th shifted version.

- Sampling π_k :

$$P(\pi_k|-) = \text{Beta}(\tilde{a}, \tilde{b}), \quad (11)$$

where $\tilde{a} = \sum_{n=1}^N b_{nk} + \frac{1}{K}$, $\tilde{b} = N + b - \sum_{n=1}^N b_{nk}$.

- Sampling γ_n :

$$P(\gamma_n|-) = \text{Gamma}(\tilde{c}, \tilde{d}), \quad (12)$$

where $\tilde{c} = c + \frac{1}{2}$, $\tilde{d} = d + \frac{1}{2} \|\mathbf{X}_n\|_2^2$.

- Sampling α_{nki} :

$$P(\alpha_{nki}|-) = \text{Gamma}(\tilde{e}, \tilde{f}), \quad (13)$$

where $\tilde{e} = e + \frac{1}{2}$, $\tilde{f} = f + \frac{1}{2} w_{nki}^2$.

- Sampling β_{kj} :

$$P(\beta_{kj}|-) = \text{Gamma}(\tilde{g}, \tilde{h}), \quad (14)$$

where $\tilde{g} = g + \frac{1}{2}$, $\tilde{h} = h + \frac{1}{2} d_{kj}^2$.

During the inference, we perform efficient block Gibbs sampling by sampling each dictionary element and corresponding hidden units given other parameters in each layer. Throughout the inference procedure, there is no matrix inversion since the dictionary and hidden units (factor scores) have been factorized. We explicitly utilize the dictionary element matrix $\mathbf{d}_k \in \mathbb{R}^{n'_x \times n'_y}$ instead of the zero-padded dictionary matrix, which is far smaller than the size of images. Moreover, thanks to the sparseness of the binary matrix b_{nk} and the max-pooling operation, the sparse matrix \mathbf{W}_n only requires modest computer memory. In practice we typically do not need a large number of Gibbs iterations to achieve useful results (although, from a pure statistics standpoint, rigorous convergence for the full posterior distribution has certainly not been achieved).

SUPPLEMENTARY MATERIAL: BATCH VB UPDATE EQUATIONS

- For b_{nk}^l we have $q(b_{nk}^l) = \text{Bernoulli}(b_{nk}^l; \rho_{nk}^l)$ where ρ_{nk}^l is the possibility of $b_{nk}^l = 1$. We consider the following two conditions:

$$\begin{aligned} \ln[q(b_{nk}^l = 1)] &\propto \zeta_1 = \langle \ln \pi_k^l \rangle - \frac{\langle \gamma_n^l \rangle}{2} (-2 \langle \|\mathbf{X}_{nk}^l \odot (\mathbf{d}_k^l * \mathbf{W}_{nk}^l)\|_2^2 \rangle + \langle \|\mathbf{W}_{nk}^l\|_2^2 \rangle \langle \|\mathbf{d}_k^l\|_2^2 \rangle) \\ \ln[q(b_{nk}^l = 0)] &\propto \zeta_2 = \langle \ln(1 - \pi_k^l) \rangle \end{aligned} \quad (15)$$

where $\mathbf{X}_{nk}^l = \mathbf{X}_{-n}^l + \mathbf{W}_{nk}^l * \mathbf{d}_k^l$, $\mathbf{X}_{-n}^l = \mathbf{X}_n^l - \sum_{k=1}^{K^l} b_{nk}^l \mathbf{W}^l * \mathbf{d}_k^l$ and b_{nk}^l is the most recent estimation, $\langle \ln \pi_k^l \rangle = \Psi(\frac{1}{K^l} + \sum_{n=1}^N \langle b_{nk}^l \rangle) - \Psi(\frac{1}{K^l} + a + N)$, $\langle \ln(1 - \pi_k^l) \rangle = \Psi(a + N - \sum_{n=1}^N \langle b_{nk}^l \rangle) - \Psi(\frac{1}{K^l} + a + N)$, $\Psi(x) = \frac{\partial}{\partial x} \ln \Gamma(x)$ and $\Gamma(x) = \int_0^\infty d\tau \tau^{x-1} e^{-\tau}$. Therefore, we can calculate $\rho_{nk}^l = \frac{1}{1 + \exp(\zeta_0 - \zeta_1)}$.

- For π_k^l we have $q(\pi_k^l) = \text{Beta}(\pi_k^l; \tau_{k1}^l, \tau_{k2}^l)$ where

$$\tau_{k1}^l = \sum_{n=1}^N \langle b_{nk}^l \rangle + \frac{1}{K^l}, \quad \tau_{k2}^l = N + b - \sum_{n=1}^N \langle b_{nk}^l \rangle \quad (16)$$

- For \mathbf{d}_k^l we have $q(\mathbf{d}_k^l) = \mathcal{N}(\mathbf{d}_k^l; \boldsymbol{\xi}_k^l, \boldsymbol{\Lambda}_k^l)$, with

$$\begin{aligned} \boldsymbol{\Lambda}_k^l &= \mathbf{1} \odot \left(\sum_{n=1}^N \langle \gamma_n^l \rangle \langle b_{nk}^l \rangle \langle \|\mathbf{W}_{nk}^l\|_2^2 \rangle + \langle \beta_k^l \rangle \right) \\ \boldsymbol{\xi}_k^l &= \boldsymbol{\Lambda}_k^l \odot \left(\sum_{n=1}^N \langle b_{nk}^l \rangle \langle \gamma_n^l \rangle (\mathbf{X}_{-n}^l * \langle \mathbf{W}_{nk}^l \rangle + \langle \mathbf{d}_k^l \rangle \langle \|\mathbf{W}_{nk}^l\|_2^2 \rangle) \right). \end{aligned} \quad (17)$$

- For \mathbf{W}_{nk}^l we have $q(w_{nki}^l) = \mathcal{N}(w_{nki}^l; \mu_{nki}^l, \Sigma_{nki}^l)$, with

$$\begin{aligned} \Sigma_{nk}^l &= \mathbf{1} \odot (\langle \gamma_n^l \rangle \langle \|\mathbf{d}_k^l\|_2^2 \rangle \langle b_{nk}^l \rangle + \langle \boldsymbol{\alpha}_{nk}^l \rangle) \\ \boldsymbol{\mu}_{nk}^l &= \langle b_{nk}^l \rangle \langle \gamma_n^l \rangle \Sigma_{nk}^l \odot (\mathbf{X}_{-n}^l * \langle \mathbf{d}_k^l \rangle + \langle \|\mathbf{d}_k^l\|_2^2 \rangle \langle \mathbf{W}_{nk}^l \rangle), \end{aligned} \quad (18)$$

where both of Σ_{nk}^l and μ_{nk}^l have the same size with \mathbf{W}_{nk}^l .

- For γ_n^l we have $q(\gamma_n^l) = \text{Gamma}(\gamma_n^l; \lambda_{n1}^l, \lambda_{n2}^l)$, where

$$\lambda_{n1}^l = c + \frac{p^l}{2}, \quad \lambda_{n2}^l = d + \frac{1}{2} \langle \|\mathbf{X}_n^l - \sum_{k=1}^{K^l} b_{nk}^l \mathbf{d}_k^l * \mathbf{W}_{nk}^l\|_2^2 \rangle \quad (19)$$

- For β_{kj}^l we have $q(\beta_{kj}^l) = \text{Gamma}(\beta_{kj}^l; v_{kj1}^l, v_{kj2}^l)$, with

$$v_{kj1}^l = g + 1/2, \quad v_{kj2}^l = h + \frac{1}{2} \langle (d_{kj}^l)^2 \rangle \quad (20)$$

- For α_{nki}^l we have $q(\alpha_{nki}^l) = \text{Gamma}(\alpha_{nki}^l; \nu_{nki1}^l, \nu_{nki2}^l)$, with

$$\nu_{nki1}^l = e + 1/2, \quad \nu_{nki2}^l = f + \frac{1}{2} \langle (w_{nki}^l)^2 \rangle \quad (21)$$