

# SUPPLEMENTARY TO: SCALABLE GROUP LEVEL PROBABILISTIC SPARSE FACTOR ANALYSIS

*Jesper L. Hinrich<sup>\*</sup>, Søren F. V. Nielsen<sup>\*</sup>, Nicolai A. B. Riis<sup>\*</sup>, Casper T. Eriksen<sup>\*</sup>, Jacob Frøsig<sup>\*</sup>,  
Marco D. F. Kristensen<sup>\*</sup>, Mikkel N. Schmidt<sup>\*</sup>, Kristoffer H. Madsen<sup>\*,†</sup> and Morten Mørup<sup>\*</sup>*

<sup>\*</sup> DTU Compute, Technical University of Denmark, Denmark

<sup>†</sup> Danish Research Centre for Magnetic Resonance, Copenhagen University Hospital Hvidovre, Denmark

## ABSTRACT

The derived details for probabilistic sparse factor analysis model for group level analysis are given. The model was first proposed by the authors and applied to functional magnetic resonance imaging[1].

### A. PROBABILISTIC SPARSE FACTOR ANALYSIS

The model's likelihood function, for a data array of size  $V \times T \times B$  with all relevant parameters collected in  $\theta$ , can be written as,

$$\mathcal{L}(\mathbf{X}|\theta) = \prod_{b=1}^B \prod_{t=1}^T \mathcal{N}\left(\mathbf{x}_t^{(b)} | \mathbf{A}\mathbf{s}_t^{(b)} + \boldsymbol{\mu}^{(b)}, \text{diag}\left(\boldsymbol{\tau}^{(b)}\right)^{-1}\right), \quad (1)$$

in which  $\mathbf{x}_t^{(b)}$  is a vector of length  $V$ ,  $\mathbf{A}$  is a matrix of size  $V \times D$ , where  $D$  is the size of the latent space,  $\mathbf{s}_t^{(b)}$  is a vector of size  $D$  and  $\boldsymbol{\tau}^{(b)}$  is a vector of length  $V$ . The distributions of the parameters in the model are,

$$\begin{aligned} P(\mathbf{A}|\boldsymbol{\alpha}) &= \prod_{v=1}^V \mathcal{N}\left(\mathbf{a}_v | \mathbf{0}, \text{diag}(\boldsymbol{\alpha}_v)^{-1}\right) \\ P(\mathbf{S}|\boldsymbol{\gamma}) &= \prod_{b=1}^B \prod_{t=1}^T \mathcal{N}\left(\mathbf{s}_t^{(b)} | \mathbf{0}, \text{diag}(\boldsymbol{\gamma})^{-1}\right) \\ P(\boldsymbol{\mu}) &= \prod_{b=1}^B \mathcal{N}\left(\boldsymbol{\mu}^{(b)} | \mathbf{0}, \beta^{-1} \mathbf{I}_V\right) \\ P(\boldsymbol{\tau}) &= \prod_{b=1}^B \prod_{v=1}^V \mathcal{G}\left(\tau_v^{(b)} | a_\tau, b_{\tau_v^{(b)}}\right) \\ P(\boldsymbol{\gamma}) &= \prod_{d=1}^D \mathcal{G}\left(\gamma_d | a_\gamma, b_{\gamma_d}\right) \\ P(\boldsymbol{\alpha}) &= \prod_{v=1}^V \prod_{d=1}^D \mathcal{G}\left(\alpha_{vd} | a_\alpha, b_{\alpha_{vd}}\right). \end{aligned}$$

Finding the parameters  $\theta$  from observed data  $\mathbf{X}$  can be done by inferring the posterior distribution,

$$P(\theta|\mathbf{X}) = \mathcal{L}(\mathbf{X}|\theta)P(\theta)/P(\mathbf{X}).$$

Unfortunately an exact inference is unfeasible for all but the simplest problems. Therefore an approximate solution is sought. While there are numerous ways to tackle this, we use variational Bayesian (VB) inference as used by [2] for VB principal components analysis. We use a mean field approximation, and find the following variational distribution to approximate the posterior.

$$\begin{aligned} Q(\mathbf{A}) &= \prod_{v=1}^V \mathcal{N}(\mathbf{a}_v | \boldsymbol{\mu}_{\mathbf{A}_v}, \boldsymbol{\Sigma}_{\mathbf{A}}^{(v)}) \\ Q(\mathbf{S}) &= \prod_{b,t=1}^{B,T} \mathcal{N}(\mathbf{s}_t^{(b)} | \boldsymbol{\mu}_{\mathbf{S}_t}^{(b)}, \boldsymbol{\Sigma}_{\mathbf{S}}^{(b)}) \\ Q(\boldsymbol{\mu}) &= \prod_{b=1}^B \mathcal{N}\left(\boldsymbol{\mu}^{(b)} | \boldsymbol{\mu}_{\boldsymbol{\mu}}^{(b)}, \boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{(b)}\right) \\ Q(\boldsymbol{\tau}) &= \prod_{b,v=1}^{B,V} \mathcal{G}(\tau_v^{(b)} | \tilde{a}_\tau, \tilde{b}_{\tau_v^{(b)}}) \\ Q(\boldsymbol{\gamma}) &= \prod_{d=1}^D \mathcal{G}(\gamma_d | \tilde{a}_\gamma, \tilde{b}_{\gamma_d}) \\ Q(\boldsymbol{\alpha}) &= \prod_{d,v=1}^{D,v} \mathcal{G}(\alpha_{vd} | \tilde{a}_\alpha, \tilde{b}_{\alpha_{vd}}) \end{aligned}$$

The moments of a distribution are then found by conditioning them on all other distributions and using free-form optimization (see [2]).

### B. UPDATE RULES

The found moments of the distributions are updated cyclically using Expectation-Maximization. In each iteration, after

---

Corresponding author: sfvn@dtu.dk. This work was supported by the Lundbeck Foundation, grant no. R105-9813. The Tesla K40 GPU used was donated by the NVIDIA Corporation.

all distributions have been updated, the evidence lowerbound (ELBO) is calculated. After a number of iterations, when the relative change in ELBO is below a given threshold a set of local optimal parameters is identified.

$$\begin{aligned}
\Sigma_{\mathbf{A}}^v &= \left( \text{diag} \langle \alpha_v \rangle + \sum_{b=1}^B \langle \tau_v^{(b)} \rangle \langle \mathbf{S}^{(b)} \mathbf{S}^{(b)\top} \rangle \right)^{-1} \\
\mu_{\mathbf{A}}^v &= \Sigma_{\mathbf{A}}^v \left( \sum_{b=1}^B \langle \tau_v^{(b)} \rangle \sum_{t=1}^T \langle \mathbf{s}_t^{(b)} \rangle (x_{tv}^{(b)} - \langle \mu_v^{(b)} \rangle) \right) \\
\Sigma_{\mathbf{S}}^{(b)} &= \left( \text{diag} \langle \gamma \rangle + \langle \mathbf{A}^\top \text{diag}(\boldsymbol{\tau}^{(b)}) \mathbf{A} \rangle \right)^{-1} \\
\mu_{\mathbf{S}_t}^{(b)} &= \Sigma_{\mathbf{S}}^{(b)} \langle \mathbf{A}^\top \rangle \text{diag} \langle \boldsymbol{\tau}^{(b)} \rangle (\mathbf{x}_t^{(b)} - \langle \boldsymbol{\mu}^{(b)} \rangle) \\
\Sigma_{\boldsymbol{\mu}}^{(b)} &= (\beta \mathbf{I}_V + \text{diag} \langle \boldsymbol{\tau}^{(b)} \rangle)^{-1} \\
\mu_{\boldsymbol{\mu}}^{(b)} &= \Sigma_{\boldsymbol{\mu}}^{(b)} \text{diag} \langle \boldsymbol{\tau}^{(b)} \rangle \sum_{t=1}^T (\mathbf{x}_t^{(b)} - \langle \mathbf{A} \rangle \langle \mathbf{s}_t^{(b)} \rangle) \\
\tilde{a}_\alpha &= a_\alpha + \frac{1}{2}, \quad \tilde{b}_{\alpha_{vd}} = b_{\alpha_{vd}} + \langle a_{vd}^2 \rangle \\
\tilde{a}_\gamma &= a_\gamma + \frac{1}{2} \sum_{b=1}^B T^{(b)}, \quad \tilde{b}_{\gamma_d} = b_{\gamma_d} + \frac{1}{2} \sum_{b=1}^B \text{trace}(\langle \mathbf{s}_d \mathbf{s}_d^\top \rangle) \\
\tilde{a}_{\tau^{(b)}} &= a_\tau + \frac{T^{(b)}}{2} \\
\tilde{b}_{\tau_v^{(b)}} &= b_{\tau_v^{(b)}} + \frac{1}{2} \left[ \|\mathbf{x}_v^{(b)}\|_{\text{Fro}}^2 + T^{(b)} \langle \mu_v^{(b)2} \rangle \right. \\
&\quad - 2 \left( \langle \mathbf{a}_v \rangle \langle \mathbf{S}^{(b)} \rangle + \langle \mu_v^{(b)} \rangle \right) \mathbf{x}_v^\top \\
&\quad + 2 \langle \mathbf{a}_v \rangle \langle \mathbf{S}^{(b)} \rangle \mathbf{1}_{T^{(b)}} \langle \mu_v^{(b)} \rangle \\
&\quad \left. + \text{trace} \left( \langle \mathbf{a}_v^\top \mathbf{S}^{(b)} \mathbf{S}^{(b)\top} \mathbf{a}_v \rangle \right) \right]
\end{aligned}$$

Note  $\langle \cdot \rangle$  is the expected value under the variational distributions. Further, using the properties of the trace operator, the expected value of the expression in  $\Sigma_{\mathbf{S}}^{(b)}$  and  $\tilde{b}_{\tau_v^{(b)}}$  are determined to be,

$$\begin{aligned}
\text{trace} \left( \langle \mathbf{a}_v^\top \mathbf{S}^{(b)} \mathbf{S}^{(b)\top} \mathbf{a}_v \rangle \right) &= \text{trace} \left( \langle \mathbf{S}^{(b)} \mathbf{S}^{(b)\top} \rangle \Sigma_{\mathbf{A}}^v \right) \\
&\quad + \langle \mathbf{a}_v^\top \rangle \langle \mathbf{S}^{(b)} \mathbf{S}^{(b)\top} \rangle \langle \mathbf{a}_v \rangle
\end{aligned}$$

and

$$\begin{aligned}
\langle \mathbf{A}^\top \text{diag}(\boldsymbol{\tau}^{(b)}) \mathbf{A} \rangle &= \left( \sum_{v=1}^V \Sigma_{\mathbf{A}}^v \langle \tau_v^{(b)} \rangle \right) \\
&\quad + \langle \mathbf{A}^\top \rangle \text{diag} \langle \boldsymbol{\tau}^{(b)} \rangle \langle \mathbf{A} \rangle.
\end{aligned}$$

### C. EVIDENCE LOWERBOUND (ELBO)

The evidence lowerbound is the sum of all the expression in this section. It can be divided into to categories; 1) the expected value of the P-distributions under the Q-distributions

(i.e. substituting the moments of the P-distributions for the moments of the Q-distributions). 2) The entropy of the Q-distributions.

For each P-distribution the expected value the corresponding Q-distribution is given below,

$$\begin{aligned}
\langle \log P(\mathbf{A}|\boldsymbol{\alpha}) \rangle &= \sum_{v=1}^V -\frac{1}{2} \log(2\pi) + \frac{1}{2} \langle \log \alpha_{vd} \rangle - \frac{1}{2} \langle \alpha_{vd} \rangle \langle a_{vd}^2 \rangle \\
\langle \log P(\mathbf{S}|\boldsymbol{\gamma}) \rangle &= \sum_{b=1}^B \sum_{t=1}^{T^{(b)}} -\frac{D}{2} \log(2\pi) + \left( \frac{1}{2} \sum_{d=1}^D \langle \log \gamma_d \rangle \right) \\
&\quad - \frac{1}{2} \text{trace} \left( \text{diag} \langle \boldsymbol{\gamma} \rangle \langle \mathbf{s}_t^{(b)} \mathbf{s}_t^{(b)\top} \rangle \right) \\
\langle \log P(\boldsymbol{\mu}) \rangle &= \sum_{b=1}^B -\frac{V}{2} \log(2\pi) + \frac{V}{2} \log(\beta) - \frac{1}{2} \beta \langle \boldsymbol{\mu}_{\boldsymbol{\mu}}^{(b)\top} \boldsymbol{\mu}_{\boldsymbol{\mu}}^{(b)} \rangle \\
\langle \log P(\boldsymbol{\alpha}) \rangle &= \sum_{d=1}^D \sum_{v=1}^V -\log(\Gamma(a_\alpha)) + a_\alpha \log(b_{\alpha_{vd}}) \\
&\quad + (a_\alpha - 1) \langle \log \alpha_{vd} \rangle - b_{\alpha_{vd}} \langle \alpha_{vd} \rangle \\
\langle \log P(\boldsymbol{\gamma}) \rangle &= \sum_{d=1}^D -\log(\Gamma(a_\gamma)) + a_\gamma \log(b_{\gamma_d}) \\
&\quad + (a_\gamma - 1) \langle \log \gamma_d \rangle - b_{\gamma_d} \langle \gamma_d \rangle \\
\langle \log P(\boldsymbol{\tau}) \rangle &= \sum_{b=1}^B \sum_{v=1}^V -\log(\Gamma(a_\tau)) + a_\tau \log(b_{\tau_v^{(b)}}) \\
&\quad + (a_\tau - 1) \langle \log \tau_v^{(b)} \rangle - b_{\tau_v^{(b)}} \langle \tau_v^{(b)} \rangle
\end{aligned}$$

Finally, the entropy for each Q-distribution is,

$$\begin{aligned}
- \langle \log Q(\mathbf{A}) \rangle &= \sum_{v=1}^V \left[ \frac{1}{2} \log |\Sigma_{\mathbf{A}}^v| + \frac{D}{2} (1 + \log(2\pi)) \right] \\
- \langle \log Q(\mathbf{S}) \rangle &= \sum_{b=1}^B \sum_{t=1}^{T^{(b)}} \left[ \frac{1}{2} \log |\Sigma_{\mathbf{S}}^{(b)}| + \frac{D}{2} (1 + \log(2\pi)) \right] \\
- \langle \log Q(\boldsymbol{\mu}) \rangle &= \sum_{b=1}^B \left[ \frac{1}{2} \log |\Sigma_{\boldsymbol{\mu}}^{(b)}| + \frac{V}{2} (1 + \log(2\pi)) \right] \\
- \langle \log Q(\boldsymbol{\alpha}) \rangle &= \sum_{d=1}^D \sum_{v=1}^V \log(\Gamma(\tilde{a}_{\alpha_{vd}})) - (\tilde{a}_\alpha - 1) \psi(\tilde{a}_\alpha) \\
&\quad - \log(\tilde{b}_{\alpha_{vd}}) + \tilde{a}_\alpha \\
- \langle \log Q(\boldsymbol{\gamma}) \rangle &= \sum_{d=1}^D \log(\Gamma(\tilde{a}_\gamma)) - (\tilde{a}_\gamma - 1) \psi(\tilde{a}_\gamma) - \log(\tilde{b}_{\gamma_d}) + \tilde{a}_\gamma \\
- \langle \log Q(\boldsymbol{\tau}) \rangle &= \sum_{b=1}^B \sum_{v=1}^V \log(\Gamma(\tilde{a}_{\tau^{(b)}})) - (\tilde{a}_{\tau^{(b)}} - 1) \psi(\tilde{a}_{\tau^{(b)}}) \\
&\quad - \log(\tilde{b}_{\tau_v^{(b)}}) + \tilde{a}_{\tau^{(b)}}
\end{aligned}$$

## D. IMPLEMENTATION

A MATLAB implementation is available<sup>1</sup>. The implementation is limited to the use of a single GPU card, as well as analysis with  $T^{(i)} = T^{(j)}, \forall i, j$ . These limitations were deemed acceptable for the work in [1], as subjects with differing timesteps are not widespread in the field.

## E. REFERENCES

- [1] Jesper L. Hinrich, Søren F.V. Nielsen, Nicolai A. B. Riis, Casper T. Eriksen, Jacob Frøsig, Marco D. F. Kristensen, Mikkel N. Schmidt, Kristoffer Hougaard Madsen, and Morten Mørup. Scalable group level probabilistic sparse factor analysis. In *(in review)*.
- [2] C M Bishop. Variational principal components. In *Proceedings of the 1999 the 9th International Conference on 'Artificial Neural Networks (ICANN99)*, pages 509–514. IEEE, 1999.

---

<sup>1</sup><https://...>