

Supplementary Materials for Multilingual Factor Analysis

Anonymous ACL submission

A Joint Distribution

We show the form of the joint distribution for 2 views. Concatenating our data and parameters as below, we can use Equation (3) of (Ghahramani et al., 1996) to write

$$\begin{aligned} \mathbf{m} &= \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}, \mathbf{W} = \begin{bmatrix} \mathbf{W}_x \\ \mathbf{W}_y \end{bmatrix} \\ \mathbf{\Psi} &= \begin{bmatrix} \mathbf{\Psi}_x & \mathbf{0} \\ \mathbf{0} & \mathbf{\Psi}_y \end{bmatrix}, \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix} \\ p(\mathbf{m}, \mathbf{z} | \boldsymbol{\theta}) &= \mathcal{N} \left(\begin{bmatrix} \mathbf{m} \\ \mathbf{z} \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{0} \end{bmatrix}, \boldsymbol{\Sigma}_{\mathbf{m}, \mathbf{z}} \right) \quad (1) \\ \boldsymbol{\Sigma}_{\mathbf{m}, \mathbf{z}} &= \begin{bmatrix} \mathbf{W}\mathbf{W}^\top + \mathbf{\Psi} & \mathbf{W} \\ \mathbf{W}^\top & \mathbb{I} \end{bmatrix} \end{aligned}$$

It is clear that this generalises to any number of views of any dimension, as the concatenation operation does not make any assumptions.

B Projections to Latent Space $\mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\mathbf{z}]$

We can query the joint Gaussian in 1 using rules from (Petersen et al., 2008) Sections (8.1.2, 8.1.3) and we get

$$\begin{aligned} p(\mathbf{z}|\mathbf{x}) &= \mathcal{N} \left(\mathbf{z}; \mathbf{W}_x^\top \boldsymbol{\Sigma}_x^{-1} \tilde{\mathbf{x}}, \mathbb{I} - \mathbf{W}_x^\top \boldsymbol{\Sigma}_x^{-1} \mathbf{W}_x \right) \\ \mathbb{E}[\mathbf{z}|\mathbf{x}] &= \mathbf{W}_x^\top \boldsymbol{\Sigma}_x^{-1} \tilde{\mathbf{x}} \end{aligned}$$

C Derivation for the Marginal Likelihood

We want to compute $p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta})$ so that we can then learn the parameters $\boldsymbol{\theta} = \{\boldsymbol{\theta}_x, \boldsymbol{\theta}_y\}$, $\boldsymbol{\theta}_i = \{\boldsymbol{\mu}_i, \mathbf{W}_i, \mathbf{\Psi}_i\}$ by maximising the marginal likelihood as is done in Factor Analysis.

From the joint $p(\mathbf{m}, \mathbf{z} | \boldsymbol{\theta})$, again using rules from (Petersen et al., 2008) Sections (8.1.2) we get

$$\begin{aligned} p(\mathbf{m} | \boldsymbol{\theta}) &= p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) \\ &= \mathcal{N} \left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \mathbf{W}\mathbf{W}^\top + \mathbf{\Psi} \right) \end{aligned}$$

For the case of two views, the joint probability can be factored as

$$\begin{aligned} p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) &= p(\mathbf{x} | \boldsymbol{\theta}_x) p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) \\ p(\mathbf{x} | \boldsymbol{\theta}_x) &= \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x) \\ p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{y}; \mathbf{W}_y \mathbf{W}_x^\top \boldsymbol{\Sigma}_x^{-1} \tilde{\mathbf{x}} + \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_{\mathbf{y}|\mathbf{x}}) \\ &= \mathcal{N}(\mathbf{y}; \mathbf{W}_y E[\mathbf{z} | \mathbf{x}] + \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_{\mathbf{y}|\mathbf{x}}), \end{aligned}$$

where

$$\begin{aligned} \boldsymbol{\Sigma}_x &= \mathbf{W}_x \mathbf{W}_x^\top + \mathbf{\Psi}_x \\ \boldsymbol{\Sigma}_{\mathbf{y}|\mathbf{x}} &= \boldsymbol{\Sigma}_y - \mathbf{W}_y \mathbf{W}_x^\top \boldsymbol{\Sigma}_x^{-1} \mathbf{W}_x \mathbf{W}_y^\top \end{aligned}$$

D Scaled Reconstruction Errors

$$\log p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) = \log p^*(\mathbf{x} | \boldsymbol{\theta}_x) + \log p^*(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) + C$$

$$C = -\frac{1}{2} (\log |2\pi \boldsymbol{\Sigma}_{\mathbf{y}|\mathbf{x}}| + \log |2\pi \boldsymbol{\Sigma}_x|)$$

$$\log p^*(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) = -\frac{1}{2} \|\tilde{\mathbf{y}} - \mathbf{W}_y E[\mathbf{z} | \mathbf{x}]\|_{\boldsymbol{\Sigma}_{\mathbf{y}|\mathbf{x}}}^2 \quad (2)$$

$$\begin{aligned} \log p^*(\mathbf{x} | \boldsymbol{\theta}_x) &= -\frac{1}{2} \|\mathbf{x} - \boldsymbol{\mu}_x\|_{\boldsymbol{\Sigma}_x}^2 \\ &= -\frac{1}{2} \|\boldsymbol{\Sigma}_x^{-\frac{1}{2}} \tilde{\mathbf{x}}\|^2 \end{aligned}$$

Setting $\mathbf{A} = \mathbf{\Psi}_x \boldsymbol{\Sigma}_x^{-1} \mathbf{\Psi}_x$, we can re-parametrise as

$$\begin{aligned} \log p^*(\mathbf{x} | \boldsymbol{\theta}_x) &= -\frac{1}{2} \|\mathbf{\Psi}_x \boldsymbol{\Sigma}_x^{-1} \tilde{\mathbf{x}}\|_{\mathbf{A}}^2 \\ &= -\frac{1}{2} \|(\boldsymbol{\Sigma}_x - \mathbf{W}_x \mathbf{W}_x^\top) \boldsymbol{\Sigma}_x^{-1} \tilde{\mathbf{x}}\|_{\mathbf{A}}^2 \\ &= -\frac{1}{2} \|\tilde{\mathbf{x}} - \mathbf{W}_x \mathbf{W}_x^\top \boldsymbol{\Sigma}_x^{-1} \tilde{\mathbf{x}}\|_{\mathbf{A}}^2 \\ &= -\frac{1}{2} \|\tilde{\mathbf{x}} - \mathbf{W}_x E[\mathbf{z} | \mathbf{x}]\|_{\mathbf{A}}^2 \end{aligned}$$

E Expectation Maximisation for MBFA

Define

$$\tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{x}_1 - \boldsymbol{\mu}_1 \\ \vdots \\ \mathbf{x}_v - \boldsymbol{\mu}_1 \end{bmatrix}, \mathbf{W} = \begin{bmatrix} \mathbf{W}_1 \\ \vdots \\ \mathbf{W}_v \end{bmatrix}$$

$$\boldsymbol{\Psi} = \begin{bmatrix} \boldsymbol{\Psi}_1 & & 0 \\ & \ddots & \\ 0 & & \boldsymbol{\Psi}_v \end{bmatrix} = \text{Bdiag}(\boldsymbol{\Psi}_1, \dots, \boldsymbol{\Psi}_v)$$

Hence

$$p(\tilde{\mathbf{x}}|\mathbf{z}; \boldsymbol{\Psi}, \mathbf{W}) = \mathcal{N}(\tilde{\mathbf{x}}|\mathbf{W}\mathbf{z}, \boldsymbol{\Psi})$$

This follows the same form as regular factor analysis, but with a block-diagonal constraint on $\boldsymbol{\Psi}$. Thus by Equations (5) and (6) of (Ghahramani et al., 1996), we apply EM as follows.

E-Step: Compute $\mathbb{E}[\mathbf{z}|\mathbf{x}]$ and $\mathbb{E}[\mathbf{z}\mathbf{z}^\top|\mathbf{x}]$ given the parameters $\boldsymbol{\theta}_t = \{\mathbf{W}_t, \boldsymbol{\Psi}_t\}$.

$$\begin{aligned} \mathbb{E}[\mathbf{z}^{(i)}|\tilde{\mathbf{x}}^{(i)}] &= \mathbf{B}_t \tilde{\mathbf{x}}^{(i)} \\ \mathbb{E}[\mathbf{z}^{(i)}\mathbf{z}^{(i)\top}|\tilde{\mathbf{x}}^{(i)}] &= \mathbb{I} - \mathbf{B}_t \mathbf{W}_t + \mathbf{B}_t \tilde{\mathbf{x}}^{(i)} \tilde{\mathbf{x}}^{(i)\top} \mathbf{B}_t^\top \\ &= \mathbf{M}_t + \mathbf{B}_t \tilde{\mathbf{x}}^{(i)} \tilde{\mathbf{x}}^{(i)\top} \mathbf{B}_t^\top \end{aligned} \quad (3)$$

where

$$\begin{aligned} \mathbf{M}_t &= \left(\mathbb{I} + \mathbf{W}_t^\top \boldsymbol{\Psi}_t^{-1} \mathbf{W}_t \right)^{-1} \\ \mathbf{B}_t &= \mathbf{W}_t^\top (\boldsymbol{\Psi}_t + \mathbf{W}_t \mathbf{W}_t^\top)^{-1} \\ &= \mathbf{M}_t \mathbf{W}_t^\top \boldsymbol{\Psi}_t^{-1}. \end{aligned} \quad (4)$$

Equation 3 is obtained by applying the Woodbury identity, and Equation 4 by applying the closely related push-through identity, as found in Section 3.2 of (Petersen et al., 2008).

M-Step: Update parameters $\boldsymbol{\theta}_{t+1} = \{\mathbf{W}_{t+1}, \boldsymbol{\Psi}_{t+1}\}$.

Define

$$\mathbf{S} = \frac{1}{m} \sum_{i=1}^m \tilde{\mathbf{x}}^{(i)} \tilde{\mathbf{x}}^{(i)\top}$$

By first observing

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \tilde{\mathbf{x}}^{(i)} \mathbb{E}[\mathbf{z}^{(i)}|\tilde{\mathbf{x}}^{(i)}]^\top &= \mathbf{S} \mathbf{B}_t^\top \\ \frac{1}{m} \sum_{j=1}^m \mathbb{E}[\mathbf{z}^{(j)}\mathbf{z}^{(j)\top}|\tilde{\mathbf{x}}^{(j)}] &= \mathbf{M}_t + \mathbf{B}_t \mathbf{S} \mathbf{B}_t^\top, \end{aligned}$$

update the parameters as follows.

$$\begin{aligned} \mathbf{W}_{t+1} &= \mathbf{S} \mathbf{B}_t^\top \left(\mathbb{I} - \mathbf{B}_t \mathbf{W}_t + \mathbf{B}_t \mathbf{S} \mathbf{B}_t^\top \right)^{-1} \\ &= \mathbf{S} \mathbf{B}_t^\top \left(\mathbf{M}_t + \mathbf{B}_t \mathbf{S} \mathbf{B}_t^\top \right)^{-1} \\ \tilde{\boldsymbol{\Psi}}_{t+1} &= \frac{1}{m} \left(\sum_{i=1}^m \tilde{\mathbf{x}}^{(i)} \tilde{\mathbf{x}}^{(i)\top} - \mathbf{W}_{t+1} \mathbb{E}[\mathbf{z}^{(i)}|\tilde{\mathbf{x}}^{(i)}] \tilde{\mathbf{x}}^{(i)\top} \right) \\ &= \mathbf{S} - \frac{1}{m} \sum_{i=1}^m \mathbf{W}_{t+1} \mathbf{B}_t \tilde{\mathbf{x}}^{(i)} \tilde{\mathbf{x}}^{(i)\top} \\ &= \mathbf{S} - \mathbf{W}_{t+1} \mathbf{B}_t \mathbf{S} \\ &= \mathbf{S} - \mathbf{S} \mathbf{B}_t^\top \mathbf{W}_{t+1}^\top \end{aligned}$$

Imposing the block diagonal constraint,

$$\boldsymbol{\Psi}_{t+1} = \text{Bdiag} \left((\tilde{\boldsymbol{\Psi}}_{t+1})_{11}, \dots, (\tilde{\boldsymbol{\Psi}}_{t+1})_{vv} \right)$$

where $(\tilde{\boldsymbol{\Psi}})_{ii} = \boldsymbol{\Psi}_i$.

F Independence to Noise in Direct Methods

We are maximising the following quantity with respect to $\boldsymbol{\theta} = \{\mathbf{W}, \boldsymbol{\mu}, \boldsymbol{\Psi}\}$

$$\begin{aligned} p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) &= \prod_i p(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}) \\ &= \prod_i \mathcal{N}(\mathbf{y}^{(i)}; \mathbf{W}\mathbf{x}^{(i)} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \\ \log p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) &= -\frac{1}{2} \left(\sum_i \|\mathbf{y}^{(i)} - \mathbf{W}\mathbf{x}^{(i)}\|_{\boldsymbol{\Psi}}^2 - C \right) \\ \frac{\partial \log p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta})}{\partial \mathbf{W}} &\propto \left(\sum_i \boldsymbol{\Psi}^{-1} (\mathbf{y}^{(i)} - \mathbf{W}\mathbf{x}^{(i)}) \mathbf{x}^{(i)\top} \right) \\ &\propto \boldsymbol{\Psi}^{-1} \left(\sum_i \mathbf{y}^{(i)} \mathbf{x}^{(i)\top} - \mathbf{W} \sum_i \mathbf{x}^{(i)} \mathbf{x}^{(i)\top} \right) \end{aligned}$$

The maximum likelihood is achieved when

$$\frac{\partial \log p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta})}{\partial \mathbf{W}} = \mathbf{0},$$

and since $\boldsymbol{\Psi}^{-1}$ has an inverse (namely $\boldsymbol{\Psi}$), this means that

$$\mathbf{W} \sum_i \mathbf{x}^{(i)} \mathbf{x}^{(i)\top} = \sum_i \mathbf{y}^{(i)} \mathbf{x}^{(i)\top}$$

It is clear from here that the MLE of \mathbf{W} does not depend on $\boldsymbol{\Psi}$, thus we can conclude that adding a noise parameter to this directed linear model has no effect on its predictions.

Table 1: Precision @1 between MBFA fitted for 1K iterations and MBFA fitted for 20K iterations.

Method	EN-IT	IT-EN	EN-FR	FR-EN	IT-FR	FR-IT
MBFA-1K	71.9	73.3	76.7	78.2	82.4	77.5
MBFA-20K	71.9	73.4	76.7	78.1	82.6	77.5
MBFA-1K+CSLS	77.5	77.6	81.9	82.0	86.8	82.1
MBFA-20K+CSLS	77.4	77.7	81.9	82.1	86.8	81.9

G Learning curve of EM

Figure 1 shows the negative log-likelihood of the three language model over the first 5,000 iterations. The precision of the learned model is very close when evaluated at iteration 1,000 and at iteration 20,000 as seen in Table 1. This suggests that the model need not be trained to full convergence to work well.

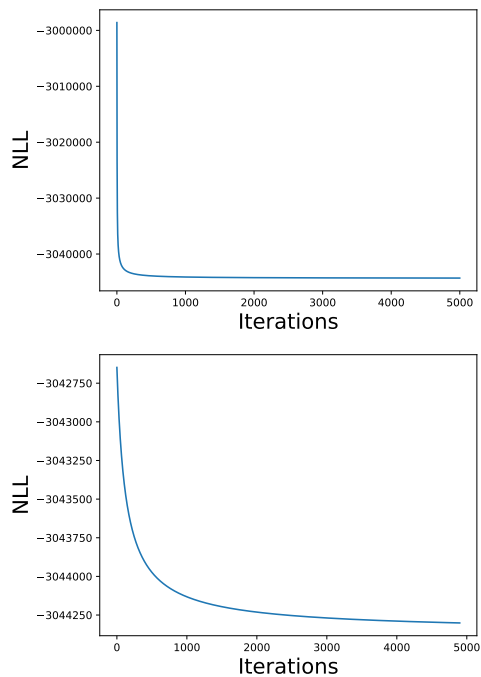


Figure 1: Training curve of EM algorithm over the first 5,000 iterations. It is clear that the procedure quickly finds a good approximation to the optimal parameters and then slowly converges to the real optimum. Top picture shows the entire training curve, while the bottom picture starts from iteration 100.

References

- Zoubin Ghahramani, Geoffrey E Hinton, et al. 1996. The em algorithm for mixtures of factor analyzers. Technical report, Technical Report CRG-TR-96-1, University of Toronto.

Kaare Brandt Petersen, Michael Syskind Pedersen, et al. 2008. The matrix cookbook. *Technical University of Denmark*, 7(15):510.