

# Bayesian Factor Analysis

Carlos Sevilla-Salcedo\*

February 2021

In this document I will extend and explain the bayesian formulation of two well known factor analysis models, Principal Component Analysis (PCA) and Canonical Correlation Analysis (CCA). The main objective of this document is to develop the equations and justify the complete distributions of the random variables of both these models.

One of the most used extensions of the probabilistic formulation is the bayesian approach. This does not only provide prior information that include previous information to the model, but also develops the posterior distribution of the observed data to sample from that distributions. This implies that we are not only estimating the value of each parameter but also the distribution that describes each one of them. The process to estimate the distribution of these variables is known as bayesian variational inference.

First of all, in order to work with bayesian variational inference we need to know the posterior distribution of the variables,  $p(\Theta|\mathbf{X})$ , where  $\Theta$  is the group of all the variables. However, this posterior distribution is not tractable. Nevertheless, we can approximate it by a distribution  $q(\Theta)$ , being  $q(\Theta)$  a group of treatable variables.

In order to adjust  $q(\Theta)$ , we need to minimise the following equation:

$$KL(q||p) = \int q(\Theta) \ln\left(\frac{q(\Theta)}{p(\Theta|\mathbf{X})}\right) d\Theta, \quad (1)$$

where KL symbolizes the Kullback-Leibler divergence, which measures the difference between two distributions. This implies that we minimise the distance between these two distributions. If we develop this equation, we have that

$$\begin{aligned} KL(q||p) &= \int q(\Theta) \ln\left(\frac{q(\Theta)}{\frac{p(\mathbf{X},\Theta)}{p(\mathbf{X})}}\right) d\Theta \\ &= \int q(\Theta) \left( \ln\left(\frac{q(\Theta)}{p(\mathbf{X},\Theta)}\right) + \ln(p(\mathbf{X})) \right) d\Theta \\ &= \int q(\Theta) \ln\left(\frac{q(\Theta)}{p(\mathbf{X},\Theta)}\right) d\Theta + \int q(\Theta) \ln(p(\mathbf{X})) d\Theta \\ &= \int q(\Theta) \ln\left(\frac{q(\Theta)}{p(\mathbf{X},\Theta)}\right) d\Theta + \ln(p(\mathbf{X})). \end{aligned} \quad (2)$$

This way, we can rearrange the previous equation to get that

$$L(q) = \ln(p(\mathbf{X})) - KL(q||p) \leq \ln(p(\mathbf{X})), \quad (3)$$

---

\*Corresponding author. Email address: sevisal@tsc.uc3m.es

having that

$$L(q) = - \int q(\Theta) \ln \left( \frac{q(\Theta)}{p(\mathbf{X}, \Theta)} \right) d\Theta \quad (4)$$

is a lower bound of  $\ln(p(\mathbf{X}))$ . For this reason, maximising the bound implies minimising  $KL(q||p)$ , keeping in mind  $L(q)$  will have a maximum value when  $p = q$ . This simple trick allows us to approximate the original variable distribution to another that can be calculated, e.g. using the mean field method.

### The mean field method

The mean field method raises as an approximation to the posterior of the model variables previously explained. Specifically, it considers that one can approximate the posterior distribution of all the model variables  $q(\Theta)$  by factorising over all the variables

$$q(\Theta) = \prod_i q(\Theta_i) = \prod_i q_i. \quad (5)$$

Hence, we can use the lower bound to determine the distribution  $q(\Theta)$ , so that it maximizes  $L(q)$ . To do so, let's substitute Equation (5) in (4)

$$\begin{aligned} L(q_j) &= \int q(\Theta) \ln \left( \frac{p(\mathbf{X}, \Theta)}{q(\Theta)} \right) d\Theta = \int \prod_i q_i \left[ \ln(p(\mathbf{X}, \Theta)) - \sum_i \ln(q_i) \right] d\Theta \\ &= \int \prod_i q_i \ln(p(\mathbf{X}, \Theta)) d\Theta - \int \prod_i q_i \sum_i \ln(q_i) d\Theta \\ &= \int q_j \prod_{i \neq j} q_i \ln(p(\mathbf{X}, \Theta)) d\Theta - \int q_j \prod_{i \neq j} q_i \left( \ln(q_j) + \sum_{i \neq j} \ln(q_i) \right) d\Theta \\ &= \int q_j \prod_{i \neq j} q_i \ln(p(\mathbf{X}, \Theta)) d\Theta - \int q_j \prod_{i \neq j} q_i \sum_{i \neq j} \ln(q_i) d\Theta - \int q_j \prod_{i \neq j} q_i \ln(q_j) d\Theta \\ &= \int q_j \left[ \int \prod_{i \neq j} q_i \ln(p(\mathbf{X}, \Theta)) d\Theta_i \right] d\Theta_j - \int q_j \ln(q_j) d\Theta_j + \text{const} \\ &= \int q_j \ln(f_j) d\Theta_j - \int q_j \ln(q_j) d\Theta_j + \text{const} \end{aligned} \quad (6)$$

where

$$\ln(f_j) = \mathbb{E}_{-q_j} [\ln(p(\mathbf{X}, \Theta))] + \text{const} \quad (7)$$

and  $-q_j$  means that we calculate this on all the variables except the  $j$ -th variable. It can be now seen that Equation (6) is a negative KL between  $q_j(\Theta_j)$  and  $f_j$ , thus to maximise  $L(q_j)$  we need to minimise  $KL(q_j||f_j)$ . Therefore, we have that the optimum solution has the following expression:

$$\ln(q_j^*) = \mathbb{E}_{-q_j} [\ln(p(\mathbf{X}, \Theta))] + \text{const}. \quad (8)$$

This constitutes the basis of the variational inference and is the approximation we will use in the following approaches.

# 1 Bayesian PCA

Here we explain the bayesian extension of the PPCA first presented in [Bishop \(1999\)](#). To do so, we adapt the probabilistic version by including a prior over the parameters of the model  $(\mathbf{W}, \sigma^2)$ . We will firstly define what is the generative model, including the variable distributions and the graphic model, to later present the variational inference result obtained.

## Generative model

As we have seen before, PCA is a FA algorithm that combines the information of a set of observed data and determines a latent space. The particularity of this latent space is that its dimensions are orthogonal to one another. Keeping this in mind, we can start by defining the distribution over the observations marginalizing with bayes rules

$$p(\mathbf{X}) \sim \int p(\mathbf{X} | \mathbf{W}, \tau) p(\mathbf{W}, \tau | \mathbf{X}) d\mathbf{W} d\tau \quad (9)$$

where we have defined a new noise variable  $\tau = \frac{1}{\sigma^2}$ . This way, we can define our generative model presented in Figure 1, which includes the different variables of the model as well as the relation between them.

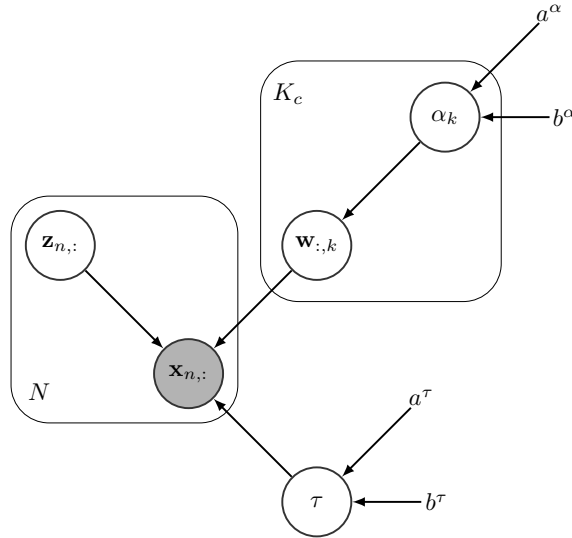


Figure 1: Plate diagram for the bayesian PCA graphical model. Gray circles denote observed variables, white circles unobserved random variables. Nodes without a circle correspond to the hyperparameters.

Furthermore, we can define the distribution of all the model variables included in the graphic model

$$\mathbf{z}_{n,:} \sim \mathcal{N}(0, \mathbf{I}_{K_c}) \quad (10)$$

$$\mathbf{w}_{:,k} \sim \mathcal{N}(0, \mathbf{I}_{K_c}) \quad (11)$$

$$\mathbf{x}_{n,:} | \mathbf{z}_{n,:} \sim \mathcal{N}(\mathbf{z}_{n,:} \mathbf{W}^T, \tau^{-1} \mathbf{I}_D) \quad (12)$$

$$\tau \sim \Gamma(a^\tau, b^\tau) \quad (13)$$

where we are assuming the observed data to be independent, which implies that the random noise modeled is the same for every sample. Note that the distribution of the observed data,  $\mathbf{x}_{n,:}$ , connects the different random variables (r.v.).

## 1.1 Variational inference

Making use of the variational inference definition previously stated, the problem can be defined to apply those techniques to design the model. First of all, considering Figure 1 we can extract the probability distribution of the model

$$p(\Theta | \mathbf{X}) \approx q(\mathbf{W})q(\tau) \prod_{n=1}^N q(\mathbf{z}_{n,:}) \quad (14)$$

This way, we can determine the optimum distribution of the different variables, applying Equation (8) and then calculate the lower bound defined in Equation (6).

In this section the different distributions specified in Equation (14) will be developed making use of the variational inference.

### 1.1.1 Distribution of $\mathbf{W}$

$$\ln(q^*(\mathbf{W})) = \mathbb{E}_{\mathbf{Z},\tau}[\ln(p(\mathbf{X}, \mathbf{W}, \mathbf{Z}, \tau))] = \mathbb{E}_{\mathbf{Z},\tau}[\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau))] + \mathbb{E}[\ln(p(\mathbf{W}))] + \text{const} \quad (15)$$

We can now evaluate both terms independently and then sum the results:

$$\begin{aligned} \ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau)) &\approx \sum_{n=1}^N \ln\left(\mathcal{N}\left(\mathbf{z}_{n,:} | \mathbf{W}^T, (\tau)^{-1}I\right)\right) + \text{const} \\ &= \sum_{n=1}^N \left( -\frac{1}{2} \ln\left|(\tau)^{-1}I\right| - \frac{1}{2} (\mathbf{x}_{n,:} - \mathbf{z}_{n,:} \mathbf{W}^T) \tau (\mathbf{x}_{n,:} - \mathbf{z}_{n,:} \mathbf{W}^T)^T \right) + \text{const} \\ &= -\frac{\tau}{2} \sum_{n=1}^N (\mathbf{x}_{n,:} \mathbf{x}_{n,:}^T - 2 \mathbf{x}_{n,:} \mathbf{W} \mathbf{z}_{n,:}^T + \mathbf{z}_{n,:} \mathbf{W}^T \mathbf{W} \mathbf{z}_{n,:}^T) + \text{const} \\ &= \sum_{n=1}^N \left( \tau \mathbf{x}_{n,:} \mathbf{W} \mathbf{z}_{n,:}^T - \frac{\tau}{2} \mathbf{z}_{n,:} \mathbf{W}^T \mathbf{W} \mathbf{z}_{n,:}^T \right) + \text{const} \\ &= \sum_{n=1}^N \left( \tau \mathbf{x}_{n,:} \mathbf{W} \mathbf{z}_{n,:}^T - \frac{\tau}{2} \sum_{d=1}^D (\mathbf{z}_{n,:} \mathbf{w}_{d,:}^T \mathbf{w}_{d,:} \mathbf{z}_{n,:}^T) \right) + \text{const} \\ &= \sum_{n=1}^N \left( \tau \mathbf{x}_{n,:} \mathbf{W} \mathbf{z}_{n,:}^T - \frac{\tau}{2} \sum_{d=1}^D (\mathbf{w}_{d,:} \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \mathbf{w}_{d,:}^T) \right) + \text{const} \end{aligned} \quad (16)$$

and applying the expectation we have that

$$\mathbb{E}_{\mathbf{Z},\tau}[\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau))] = \sum_{n=1}^N \left( \langle \tau \rangle \mathbf{x}_{n,:} \mathbf{W} \langle \mathbf{z}_{n,:} \rangle - \frac{\langle \tau \rangle}{2} \sum_{d=1}^D (\mathbf{w}_{d,:} \langle \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \rangle \mathbf{w}_{d,:}^T) \right) \quad (17)$$

Regarding the second, we have that

$$\mathbb{E}[\ln(p(\mathbf{w}_{d,:}))] = \ln(p(\mathbf{w}_{d,:})) = \sum_{d=1}^D \left( -\frac{1}{2} \mathbf{w}_{d,:} \mathbf{w}_{d,:}^T \right) + \text{const} \quad (18)$$

Finally, joining Equations (17) and (18), we obtain the mean field approximation

$$\begin{aligned}
\ln(q^*(\mathbf{W})) &= \sum_{d=1}^D \left( -\frac{1}{2} \mathbf{w}_{d,:} \mathbf{w}_{d,:}^T \right) + \sum_{n=1}^N \left( \langle \tau \rangle \mathbf{x}_{n,:} \mathbf{W} \langle \mathbf{z}_{n,:} \rangle - \frac{\langle \tau \rangle}{2} \sum_{d=1}^D (\mathbf{w}_{d,:} \langle \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \rangle \mathbf{w}_{d,:}^T) \right) + \text{const} \\
&= \sum_{d=1}^D \left( -\frac{1}{2} \mathbf{w}_{d,:} \mathbf{w}_{d,:}^T + \sum_{n=1}^N \left( \langle \tau \rangle x_{n,d} \mathbf{w}_{d,:} \langle \mathbf{z}_{n,:} \rangle - \frac{\langle \tau \rangle}{2} \mathbf{w}_{d,:} \langle \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \rangle \mathbf{w}_{d,:}^T \right) \right) + \text{const} \\
&= \sum_{d=1}^D \left( -\frac{1}{2} \mathbf{w}_{d,:} (\mathbf{I} + \langle \tau \rangle \langle \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \rangle) \mathbf{w}_{d,:}^T + \langle \tau \rangle \mathbf{w}_{d,:} \langle \mathbf{Z}^T \rangle \mathbf{x}_{:,d}^T \right) + \text{const} \tag{19}
\end{aligned}$$

This way, comparing the results with the normal distribution, we can identify terms, extracting the following conclusions:

$$q^*(\mathbf{W}) = \prod_{d=1}^D (\mathcal{N}(\mathbf{w}_{d,:} | \mu_{\mathbf{w}_{d,:}}, \Sigma_{\mathbf{W}})) \tag{20}$$

where the variance is

$$\Sigma_{\mathbf{W}}^{-1} = \mathbf{I} + \langle \tau \rangle \langle \mathbf{Z}^T \mathbf{Z} \rangle \tag{21}$$

and the mean is

$$\mu_{\mathbf{w}_{d,:}} = \langle \tau \rangle \Sigma_{\mathbf{W}} \langle \mathbf{Z} \rangle^T \mathbf{x}_{:,d}^T \tag{22}$$

or in a matricial way

$$\langle \mathbf{W} \rangle = \langle \tau \rangle \mathbf{X}^T \langle \mathbf{Z} \rangle \Sigma_{\mathbf{W}} \tag{23}$$

### 1.1.2 Distribution of $\mathbf{Z}$

$$\ln(q^*(\mathbf{Z})) = \mathbb{E}_{\mathbf{W}, \tau} [\ln(p(\mathbf{X}, \mathbf{W}, \mathbf{Z}, \tau))] = \mathbb{E}_{\mathbf{W}, \tau} [\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau))] + \mathbb{E}[\ln(p(\mathbf{Z}))] \tag{24}$$

We will again start by calculating the first term (which is similar to equation (17))

$$\mathbb{E}_{\mathbf{W}, \tau} [\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau))] = \sum_{n=1}^N \left( \langle \tau \rangle \mathbf{x}_{n,:} \langle \mathbf{W} \rangle \mathbf{z}_{n,:}^T - \frac{\langle \tau \rangle}{2} \mathbf{z}_{n,:} \langle \mathbf{W}^T \mathbf{W} \rangle \mathbf{z}_{n,:}^T \right) + \text{const} \tag{25}$$

As the associated prior of this variable is the same as the previous one, the result of the second term is equivalent

$$\mathbb{E}[\ln(p(\mathbf{Z}))] = \ln(p(\mathbf{Z})) = \sum_{n=1}^N \left( -\frac{1}{2} \mathbf{z}_{n,:} \mathbf{z}_{n,:}^T \right) + \text{const} \tag{26}$$

Therefore, joining the different terms in Equations (25) and (26), the mean field approximation has the form

$$\ln(q^*(\mathbf{Z})) = \sum_{n=1}^N \left( -\frac{1}{2} \mathbf{z}_{n,:} \mathbf{z}_{n,:}^T + \langle \tau \rangle \mathbf{x}_{n,:} \langle \mathbf{W} \rangle \mathbf{z}_{n,:}^T - \frac{\langle \tau \rangle}{2} \mathbf{z}_{n,:} \langle \mathbf{W}^T \mathbf{W} \rangle \mathbf{z}_{n,:}^T \right) + \text{const} \tag{27}$$

Again, if we identify terms from the previous equation we get that

$$q^*(\mathbf{Z}) = \prod_{n=1}^N (\mathcal{N}(\mathbf{z}_{n,:} | \mu_{\mathbf{z}_{n,:}}, \Sigma_{\mathbf{Z}})) \quad (28)$$

where the variance is

$$\Sigma_{\mathbf{Z}}^{-1} = I + \sum_{m=1}^M (\langle \tau \rangle \langle \mathbf{W}^T \mathbf{W} \rangle) \quad (29)$$

and the mean is

$$\mu_{\mathbf{z}_{n,:}} = \sum_{m=1}^M (\langle \tau \rangle \mathbf{x}_{n,:} \langle \mathbf{W} \rangle \Sigma_{\mathbf{Z}}) \quad (30)$$

or in a matricial way

$$\langle \mathbf{Z} \rangle = \sum_{m=1}^M (\langle \tau \rangle \mathbf{X} \langle \mathbf{W} \rangle \Sigma_{\mathbf{Z}}) \quad (31)$$

### 1.1.3 Distribution of $\tau$

Let's now calculate the approximate distribution of the noise variable,  $\tau$ .

$$\ln(q^*(\tau)) = \mathbb{E}_{\mathbf{W}, \mathbf{Z}} [\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau))] + \mathbb{E}[\ln(p(\tau))] \quad (32)$$

Similarly to what was done in equation (17) we can now calculate the first term, although the terms that will be constant for the expectation will vary

$$\begin{aligned} \ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau)) &= \sum_{n=1}^N \ln \left( \mathcal{N}(\mathbf{z}_{n,:} | \mathbf{W}^T, (\tau)^{-1} I) \right) + \text{const} \\ &= \sum_{n=1}^N \sum_{d=1}^D \left( \frac{1}{2} \ln |\tau| - \frac{\tau}{2} (\mathbf{x}_{n,d} - \mathbf{z}_{n,:} \mathbf{w}_{d,:}^T)^2 \right) + \text{const} \\ &= \frac{DN}{2} \ln(\tau) - \frac{\tau}{2} \sum_{n=1}^N \sum_{d=1}^D \left( \mathbf{x}_{n,d}^2 - 2 \mathbf{w}_{d,:}^T \mathbf{z}_{n,:} \mathbf{x}_{n,d} + (\mathbf{z}_{n,:} \mathbf{w}_{d,:}^T)^2 \right) + \text{const} \\ &= \frac{DN}{2} \ln(\tau) - \frac{\tau}{2} \left( \sum_{n=1}^N \sum_{d=1}^D \mathbf{x}_{n,d}^2 - 2 \sum_{d=1}^D \mathbf{w}_{d,:}^T \mathbf{Z}^T \mathbf{x}_{:,d} + \sum_{d=1}^D \mathbf{w}_{d,:}^T \mathbf{Z}^T \mathbf{Z} \mathbf{w}_{d,:}^T \right) + \text{const} \\ &= \frac{DN}{2} \ln(\tau) - \frac{\tau}{2} \left( \sum_{n=1}^N \sum_{d=1}^D \mathbf{x}_{n,d}^2 - 2 \text{Tr}\{\mathbf{W} \mathbf{Z}^T \mathbf{X}\} + \text{Tr}\{\mathbf{W} \mathbf{Z}^T \mathbf{Z} \mathbf{W}^T\} \right) + \text{const} \end{aligned}$$

We can now calculate the expectation with respect the rest of the random variables

$$\begin{aligned} \mathbb{E}_{\mathbf{W}, \mathbf{Z}} [\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau))] &= \frac{DN}{2} \ln(\tau) - \frac{\tau}{2} \left( \sum_{n=1}^N \sum_{d=1}^D \mathbf{x}_{n,d}^2 - 2 \text{Tr}\{\langle \mathbf{W} \rangle \langle \mathbf{Z}^T \rangle \mathbf{X}\} + \text{Tr}\{\langle \mathbf{W}^T \mathbf{W} \rangle \langle \mathbf{Z}^T \mathbf{Z} \rangle\} \right) + \text{const} \end{aligned}$$

The second term can subsequently be determined as:

$$\mathbb{E}[\ln(p(\tau))] = \ln(p(\tau)) = -b_0^\tau \tau + (a_0^\tau - 1) \ln(\tau) + \text{const}$$

Joining them together, we have that:

$$\begin{aligned} \ln(q^*(\tau)) &= \frac{DN}{2} \ln(\tau) - \frac{\tau}{2} \left( \sum_{n=1}^N \sum_{d=1}^D x_{n,d}^2 - 2 \text{Tr}\{\langle \mathbf{W} \rangle \langle \mathbf{Z}^T \rangle \mathbf{X}\} + \text{Tr}\{\langle \mathbf{W}^T \mathbf{W} \rangle \langle \mathbf{Z}^T \mathbf{Z} \rangle\} \right) \\ &\quad - b_0^\tau \tau + (a_0^\tau - 1) \ln(\tau) + \text{const} \end{aligned}$$

So by identifying terms with the distribution, we would have that for this parameter:

$$q^*(\tau) = (\text{Gamma}(\tau | a^\tau, b^\tau)) \quad (33)$$

where the first parameter of the distribution is

$$a^\tau = \frac{DN}{2} + \alpha_0^\tau \quad (34)$$

and the second is

$$b^\tau = \beta_0^\tau + \frac{1}{2} \left( \sum_{n=1}^N \sum_{d=1}^D x_{n,d}^2 - 2 \text{Tr}\{\langle \mathbf{W} \rangle \langle \mathbf{Z}^T \rangle \mathbf{X}\} + \text{Tr}\{\langle \mathbf{W}^T \mathbf{W} \rangle \langle \mathbf{Z}^T \mathbf{Z} \rangle\} \right) \quad (35)$$

## 2 Bayesian PCA with ARD

### 2.1 Generative model

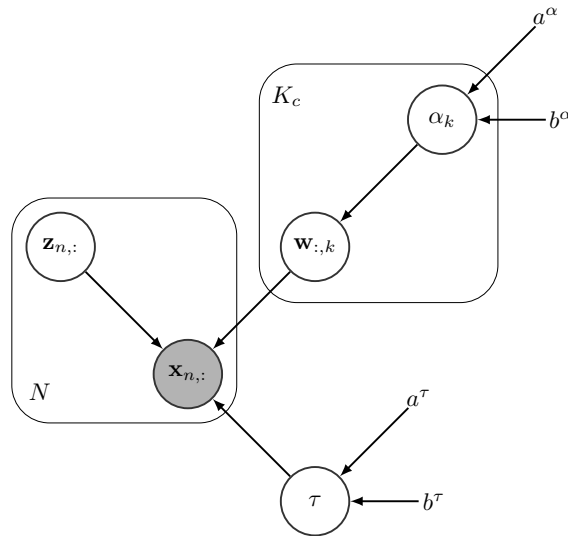


Figure 2: Plate diagram for the bayesian PCA with ARD graphical model. Gray circles denote observed variables, white circles unobserved random variables. Nodes without a circle correspond to the hyperparameters.

We can define the distribution over the observations by marginalizing using the bayes rule

$$p(\mathbf{X}) \sim \int p(\mathbf{X} | \mathbf{W}, \boldsymbol{\alpha}, \tau) p(\mathbf{W}, \boldsymbol{\alpha}, \tau | \mathbf{X}) d\mathbf{W} d\boldsymbol{\alpha} d\tau \quad (36)$$

Therefore, we can now establish the prior distributions over the different random variables in the problem

$$\mathbf{z}_{n,:} \sim \mathcal{N}(0, \mathbf{I}_{K_c}) \quad (37)$$

$$\mathbf{w}_{:,k} \sim \mathcal{N}(0, \alpha_k^{-1} \mathbf{I}_{K_c}) \quad (38)$$

$$\mathbf{x}_{n,:} | \mathbf{z}_{n,:} \sim \mathcal{N}(\mathbf{z}_{n,:} \mathbf{W}^T, \tau^{-1} \mathbf{I}_D) \quad (39)$$

$$\alpha_k \sim \Gamma(a^\alpha, b^\alpha) \quad (40)$$

$$\tau \sim \Gamma(a^\tau, b^\tau) \quad (41)$$

## 2.2 Variational inference

Making use of the variational inference definition previously stated, the problem can be defined to apply those techniques to design the model. First of all, considering Figure 2 we can extract the probability distribution of the model

$$p(\Theta | \mathbf{X}) \approx \left( q(\mathbf{W}) q(\tau) \prod_{k=1}^{K_c} q(\alpha_k) \right) \prod_{n=1}^N q(\mathbf{z}_{n,:}) \quad (42)$$

This way, we can determine the optimum distribution of the different variables, applying Equation (8) and then calculate the lower bound defined in Equation (6).

In this section the different distributions specified in Equation (42) will be developed making use of the variational inference.

## 2.3 Distribution of $\mathbf{W}$

$$\begin{aligned} \ln(q^*(\mathbf{W})) &= \mathbb{E}_{\mathbf{Z}, \boldsymbol{\alpha}, \tau} [\ln(p(\mathbf{X}, \mathbf{W}, \mathbf{Z}, \boldsymbol{\alpha}, \tau))] \\ &= \mathbb{E}_{\mathbf{Z}, \tau} [\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau))] + \mathbb{E}_{\boldsymbol{\alpha}} [\ln(p(\mathbf{W} | \boldsymbol{\alpha}))] + \text{const} \end{aligned} \quad (43)$$



We can now evaluate both terms independently and then sum the results:

$$\begin{aligned}
\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau)) &\approx \sum_{n=1}^N \ln\left(\mathcal{N}\left(\mathbf{z}_{n,:}, \mathbf{W}^T, (\tau)^{-1}I\right)\right) + \text{const} \\
&= \sum_{n=1}^N \left(-\frac{1}{2} \ln |\tau^{-1}I| - \frac{1}{2} (\mathbf{x}_{n,:} - \mathbf{z}_{n,:} \mathbf{W}^T) \tau (\mathbf{x}_{n,:} - \mathbf{z}_{n,:} \mathbf{W}^T)^T\right) + \text{const} \\
&= -\frac{\tau}{2} \sum_{n=1}^N (\mathbf{x}_{n,:} \mathbf{x}_{n,:}^T - 2 \mathbf{x}_{n,:} \mathbf{W} \mathbf{z}_{n,:}^T + \mathbf{z}_{n,:} \mathbf{W}^T \mathbf{W} \mathbf{z}_{n,:}^T) + \text{const} \\
&= \sum_{n=1}^N \left(\tau \mathbf{x}_{n,:} \mathbf{W} \mathbf{z}_{n,:}^T - \frac{\tau}{2} \mathbf{z}_{n,:} \mathbf{W}^T \mathbf{W} \mathbf{z}_{n,:}^T\right) + \text{const} \\
&= \sum_{n=1}^N \left(\tau \mathbf{x}_{n,:} \mathbf{W} \mathbf{z}_{n,:}^T - \frac{\tau}{2} \sum_{d=1}^D (\mathbf{z}_{n,:} \mathbf{w}_{d,:}^T \mathbf{w}_{d,:} \mathbf{z}_{n,:}^T)\right) + \text{const} \\
&= \sum_{n=1}^N \left(\tau \mathbf{x}_{n,:} \mathbf{W} \mathbf{z}_{n,:}^T - \frac{\tau}{2} \sum_{d=1}^D (\mathbf{w}_{d,:} \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \mathbf{w}_{d,:}^T)\right) + \text{const} \tag{44}
\end{aligned}$$

This way the expectation can be calculated as:

$$\mathbb{E}_{\mathbf{Z}, \tau}[\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau))] = \sum_{n=1}^N \left(\langle \tau \rangle \mathbf{x}_{n,:} \mathbf{W} \langle \mathbf{z}_{n,:} \rangle - \frac{\langle \tau \rangle}{2} \sum_{d=1}^D (\mathbf{w}_{d,:} \langle \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \rangle \mathbf{w}_{d,:}^T)\right) \tag{45}$$

Equivalently, the second term would be:

$$\begin{aligned}
\ln(p(\mathbf{W} | \boldsymbol{\alpha})) &= \sum_{d=1}^D \sum_{k=1}^{K_c} \ln(p(\mathbf{w}_{d,k} | \alpha_k)) \\
&= \sum_{d=1}^D \sum_{k=1}^{K_c} \left(\frac{1}{2} \ln(\alpha_k) - \frac{\alpha_k}{2} \mathbf{w}_{d,k}^2\right) + \text{const} \tag{46}
\end{aligned}$$

$$\mathbb{E}_{\boldsymbol{\alpha}}[\ln(p(\mathbf{W} | \boldsymbol{\alpha}))] = \sum_{d=1}^D \sum_{k=1}^{K_c} \left(-\frac{\langle \alpha_k \rangle}{2} \mathbf{w}_{d,k}^2\right) + \text{const} \tag{47}$$

Finally, joining Equations (45) and (47), we would have that the optimum solution for

the variable  $\mathbf{W}$  is:

$$\begin{aligned}
\ln(q^*(\mathbf{W})) &= \sum_{d=1}^D \sum_{k=1}^{K_c} \left( -\frac{\langle \alpha_k \rangle}{2} \mathbf{w}_{d,k}^2 \right) + \sum_{n=1}^N (\langle \tau \rangle \mathbf{x}_{n,:} \mathbf{W} \langle \mathbf{z}_{n,:} \rangle \\
&\quad - \frac{\langle \tau \rangle}{2} \sum_{d=1}^D (\mathbf{w}_{d,:} \langle \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \rangle \mathbf{w}_{d,:}^T)) + \text{const} \\
&= \sum_{d=1}^D \left( -\frac{1}{2} \mathbf{w}_{d,:} \text{diag}(\langle \alpha_k \rangle) \mathbf{w}_{d,:}^T + \sum_{n=1}^N (\langle \tau \rangle \mathbf{x}_{n,d} \mathbf{w}_{d,:} \langle \mathbf{z}_{n,:} \rangle \right. \\
&\quad \left. - \frac{\langle \tau \rangle}{2} \mathbf{w}_{d,:} \langle \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \rangle \mathbf{w}_{d,:}^T) \right) + \text{const} \\
&= \sum_{d=1}^D \left( -\frac{1}{2} \mathbf{w}_{d,:} (\text{diag}(\langle \alpha_k \rangle) + \langle \tau \rangle \langle \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \rangle) \mathbf{w}_{d,:}^T + \langle \tau \rangle \mathbf{w}_{d,:} \langle \mathbf{Z}^T \rangle \mathbf{x}_{:,d}^T \right) + \text{const}
\end{aligned} \tag{48}$$

This way, comparing the results with the normal distribution, we can identify terms, extracting the following conclusions:

$$q^*(\mathbf{W}) = \prod_{d=1}^D (\mathcal{N}(\mathbf{w}_{d,:} | \mu_{\mathbf{w}_{d,:}}, \Sigma_{\mathbf{W}})) \tag{49}$$

Where the variance would be:

$$\Sigma_{\mathbf{W}}^{-1} = \text{diag}(\langle \alpha \rangle) + \langle \tau \rangle \langle \mathbf{Z}^T \mathbf{Z} \rangle \tag{50}$$

And the mean could be expressed as:

$$\begin{aligned}
\mu_{\mathbf{w}_{d,:}} &= \langle \tau \rangle \Sigma_{\mathbf{W}} \langle \mathbf{Z} \rangle^T \mathbf{x}_{:,d}^T \\
\langle \mathbf{W} \rangle &= \langle \tau \rangle \mathbf{X}^T \langle \mathbf{Z} \rangle \Sigma_{\mathbf{W}}
\end{aligned} \tag{51}$$

## 2.4 Distribution of $\mathbf{Z}$

$$\ln(q^*(\mathbf{Z})) = \mathbb{E}_{\mathbf{W}, \tau} [\ln(p(\mathbf{X}, \mathbf{W}, \mathbf{Z}, \tau))] = \mathbb{E}_{\mathbf{W}, \tau} [\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau))] + \mathbb{E}[\ln(p(\mathbf{Z}))] \tag{52}$$

The first term is similar to equation (17), changing the expectation dependency:

$$\mathbb{E}_{\mathbf{W}, \tau} [\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau))] = \sum_{n=1}^N \left( \langle \tau \rangle \mathbf{x}_{n,:} \langle \mathbf{W} \rangle \mathbf{z}_{n,:}^T - \frac{\langle \tau \rangle}{2} \mathbf{z}_{n,:} \langle \mathbf{W}^T \mathbf{W} \rangle \mathbf{z}_{n,:}^T \right) + \text{const} \tag{53}$$

The second term could be calculated as:

$$\mathbb{E}[\ln(p(\mathbf{Z}))] = \ln(p(\mathbf{Z})) = \sum_{n=1}^N \left( -\frac{1}{2} \mathbf{z}_{n,:} \mathbf{z}_{n,:}^T \right) + \text{const} \tag{54}$$

Therefore, joining the different terms in Equations (53) and (54), the optimum solution has the form:

$$\ln(q^*(\mathbf{Z})) = \sum_{n=1}^N \left( -\frac{1}{2} \mathbf{z}_{n,:} \mathbf{z}_{n,:}^T + \langle \tau \rangle \mathbf{x}_{n,:} \langle \mathbf{W} \rangle \mathbf{z}_{n,:}^T - \frac{\langle \tau \rangle}{2} \mathbf{z}_{n,:} \langle \mathbf{W}^T \mathbf{W} \rangle \mathbf{z}_{n,:}^T \right) + \text{const} \tag{55}$$

This way, we can now identify terms, having that:

$$q^*(\mathbf{Z}) = \prod_{n=1}^N (\mathcal{N}(\mathbf{z}_{n,:} | \mu_{\mathbf{z}_{n,:}}, \Sigma_{\mathbf{Z}})) \quad (56)$$

Where the variance would be:

$$\Sigma_{\mathbf{Z}}^{-1} = I + \langle \tau \rangle \langle \mathbf{W}^T \mathbf{W} \rangle \quad (57)$$

And the mean could be expressed as:

$$\begin{aligned} \mu_{\mathbf{z}_{n,:}} &= \langle \tau \rangle \mathbf{x}_{n,:} \langle \mathbf{W} \rangle \Sigma_{\mathbf{Z}} \\ \langle \mathbf{Z} \rangle &= \langle \tau \rangle \mathbf{X} \langle \mathbf{W} \rangle \Sigma_{\mathbf{Z}} \end{aligned} \quad (58)$$

## 2.5 Distribution of $\alpha$

$$\ln(q^*(\alpha)) = \mathbb{E}_{\mathbf{W}}[\ln(p(\mathbf{W} | \alpha))] + \mathbb{E}[\ln(p(\alpha))] \quad (59)$$

Similarly to what we have done before, we can develop both terms independently:

$$\begin{aligned} \ln(p(\mathbf{W} | \alpha)) &= \sum_{d=1}^D \sum_{k=1}^{K_c} \left( \frac{1}{2} \ln |\alpha_k| - \frac{1}{2} \alpha_k (w_{d,k})^2 \right) + \text{const} \\ &= \sum_{k=1}^{K_c} \left( \frac{D}{2} \ln(\alpha_k) - \frac{1}{2} \alpha_k \mathbf{w}_{:,k}^T \mathbf{w}_{:,k} \right) + \text{const} \\ \mathbb{E}[\ln(p(\mathbf{W} | \alpha))] &= \sum_{k=1}^{K_c} \left( \frac{D}{2} \ln(\alpha_k) - \frac{1}{2} \alpha_k \langle \mathbf{w}_{:,k}^T \mathbf{w}_{:,k} \rangle \right) + \text{const} \end{aligned} \quad (60)$$

The second term is deployed as follows:

$$\begin{aligned} \ln(p(\alpha_k)) &= -\beta_0 \alpha_k + (\alpha_0 - 1) \ln(\alpha_k) + \text{const} \\ \mathbb{E}[\ln(p(\alpha))] &= \sum_{k=1}^{K_c} (\ln(p(\alpha_k))) = \sum_{k=1}^{K_c} (-\beta_0 \alpha_k + (\alpha_0 - 1) \ln(\alpha_k)) + \text{const} \end{aligned} \quad (61)$$

Therefore, by joining both Equations (60) and (61) together, we can now determine the value of the expectation we were looking for

$$\ln(q^*(\alpha)) = \sum_{k=1}^{K_c} \left( \frac{D}{2} \ln(\alpha_k) - \frac{1}{2} \alpha_k \langle \mathbf{w}_{:,k}^T \mathbf{w}_{:,k} \rangle - \beta_0 \alpha_k + (\alpha_0 - 1) \ln(\alpha_k) \right) + \text{const} \quad (62)$$

This way, we can now identify terms, having

$$q^*(\alpha) = \prod_{k=1}^{K_c} (\text{Gamma}(\alpha_k | a_{\alpha_k}, b_{\alpha_k})), \quad (63)$$

where the variance would be:

$$a_{\alpha_k} = \frac{D}{2} + \alpha_0 \quad (64)$$

and the mean is

$$b_{\alpha_k} = \beta_0 + \frac{1}{2} \langle \mathbf{w}_{:,k}^T \mathbf{w}_{:,k} \rangle = \beta_0 + \frac{1}{2} \langle \mathbf{W}^T \mathbf{W} \rangle_{k,k}, \quad (65)$$

where

$$\langle \mathbf{w}_{:,k}^T \mathbf{w}_{:,k} \rangle = \sum_{d=1}^D (w_{d,k})^2 = \langle \mathbf{W}^T \mathbf{W} \rangle_{k,k}.$$

## 2.6 Distribution of $\tau$

$$\ln(q^*(\tau)) = \mathbb{E}_{\mathbf{W}, \mathbf{Z}}[\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau))] + \mathbb{E}[\ln(p(\tau))] \quad (66)$$

Similarly to what was done in equation (45):

$$\begin{aligned} \ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau)) &= \sum_{n=1}^N \ln\left(\mathcal{N}(\mathbf{z}_{n,:} | \mathbf{W}^T, (\tau)^{-1} I)\right) + \text{const} \\ &= \sum_{n=1}^N \sum_{d=1}^D \left( \frac{1}{2} \ln |\tau| - \frac{\tau}{2} (x_{n,d} - \mathbf{z}_{n,:} \mathbf{w}_{d,:}^T)^2 \right) + \text{const} \\ &= \frac{DN}{2} \ln(\tau) - \frac{\tau}{2} \sum_{n=1}^N \sum_{d=1}^D \left( x_{n,d}^2 - 2 \mathbf{w}_{d,:}^T \mathbf{z}_{n,:} x_{n,d} + (\mathbf{z}_{n,:} \mathbf{w}_{d,:}^T)^2 \right) + \text{const} \\ &= \frac{DN}{2} \ln(\tau) - \frac{\tau}{2} \left( \sum_{n=1}^N \sum_{d=1}^D x_{n,d}^2 - 2 \sum_{d=1}^D \mathbf{w}_{d,:}^T \mathbf{Z}^T \mathbf{x}_{:,d} + \sum_{d=1}^D \mathbf{w}_{d,:}^T \mathbf{Z}^T \mathbf{Z} \mathbf{w}_{d,:}^T \right) + \text{const} \\ &= \frac{DN}{2} \ln(\tau) - \frac{\tau}{2} \left( \sum_{n=1}^N \sum_{d=1}^D x_{n,d}^2 - 2 \text{Tr}\{\mathbf{W} \mathbf{Z}^T \mathbf{X}\} + \text{Tr}\{\mathbf{W} \mathbf{Z}^T \mathbf{Z} \mathbf{W}^T\} \right) + \text{const} \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{\mathbf{W}, \mathbf{Z}}[\ln(p(\mathbf{X} | \mathbf{W}, \mathbf{Z}, \tau))] &= \frac{DN}{2} \ln(\tau) - \frac{\tau}{2} \left( \sum_{n=1}^N \sum_{d=1}^D x_{n,d}^2 - 2 \text{Tr}\{\langle \mathbf{W} \rangle \langle \mathbf{Z}^T \rangle \mathbf{X}\} \right. \\ &\quad \left. + \text{Tr}\{\langle \mathbf{W}^T \mathbf{W} \rangle \langle \mathbf{Z}^T \mathbf{Z} \rangle\} \right) + \text{const} \end{aligned}$$

The second term could subsequently be determined as:

$$\mathbb{E}[\ln(p(\tau))] = \ln(p(\tau)) = -\beta_0^T \tau + (\alpha_0^T - 1) \ln(\tau) + \text{const}$$

Joining them together, we have that:

$$\begin{aligned} \ln(q^*(\tau)) &= \frac{DN}{2} \ln(\tau) - \frac{\tau}{2} \left( \sum_{n=1}^N \sum_{d=1}^D x_{n,d}^2 \right. \\ &\quad \left. - 2 \text{Tr}\{\langle \mathbf{W} \rangle \langle \mathbf{Z}^T \rangle \mathbf{X}\} + \text{Tr}\{\langle \mathbf{W}^T \mathbf{W} \rangle \langle \mathbf{Z}^T \mathbf{Z} \rangle\} \right) \\ &\quad - \beta_0^T \tau + (\alpha_0^T - 1) \ln(\tau) + \text{const} \end{aligned}$$

So by identifying terms with the distribution, we would have that for this parameter:

$$q^*(\tau) = (\text{Gamma}(\tau|a_\tau, b_\tau)) \quad (67)$$

the variance would be:

$$a_\tau = \frac{DN}{2} + \alpha_0^\tau \quad (68)$$

and the mean could be expressed as:

$$b_\tau = \beta_0^\tau + \frac{1}{2} \left( \sum_{n=1}^N \sum_{d=1}^D x_{n,d}^2 - 2 \text{Tr}\{\langle \mathbf{W} \rangle \langle \mathbf{Z}^T \rangle \mathbf{X}\} + \text{Tr}\{\langle \mathbf{W}^T \mathbf{W} \rangle \langle \mathbf{Z}^T \mathbf{Z} \rangle\} \right) \quad (69)$$

### 3 Bayesian CCA or Bayesian Inter-Battery Factor Analysis (BIBFA)

#### 3.1 Generative model

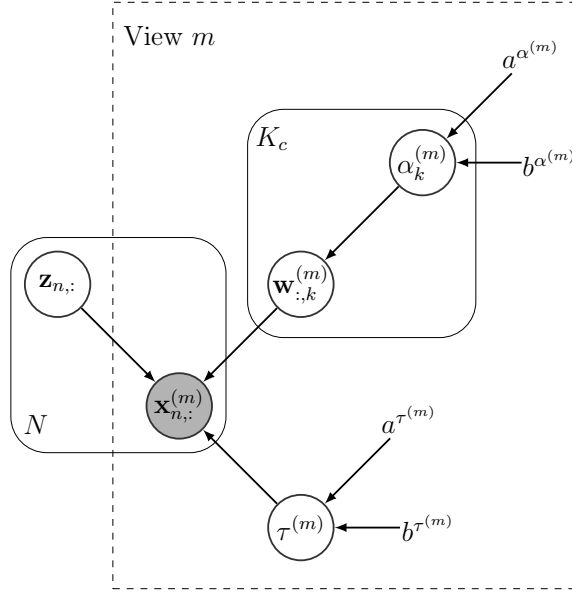


Figure 3: Graphic model of the CCA scheme for neuromarkers design.

In our case we have two views. The first one will be the input data ( $\mathbf{X}^{(1)} \in \mathbb{R}^{N \times d}$ ) and the output labels ( $\mathbf{X}^{(2)} \in \mathbb{R}^{N \times c}$ ). These two views share some latent variables ( $\mathbf{Z} \in \mathbb{R}^{N \times K}$ ) and have some specific ones ( $\mathbf{Z}^{(1)} \in \mathbb{R}^{N \times K_1}$  and  $\mathbf{Z}^{(2)} \in \mathbb{R}^{N \times K_2}$ ). This way, we have projection matrices of those shared latent variables for each view ( $\mathbf{A}^{(1)} \in \mathbb{R}^{d \times K}$  and  $\mathbf{A}^{(2)} \in \mathbb{R}^{c \times K}$ ) as well as other projection matrices for the specific latent variables ( $\mathbf{B}^{(1)} \in \mathbb{R}^{d \times K_1}$  and  $\mathbf{B}^{(2)} \in \mathbb{R}^{c \times K_2}$ ). Finally, we have the noise, which includes all we weren't able to model (the covariance matrices have the shape  $\Sigma^{(1)} \in \mathbb{R}^{d \times d}$  and  $\Sigma^{(2)} \in \mathbb{R}^{c \times c}$ ).

$$K_c = K + K_1 + K_2 \quad (70)$$

$$\left. \begin{aligned}
& \left. \begin{aligned} \mathbf{X}^{(1)} &\in \mathbb{R}^{N \times d} \\ \mathbf{X}^{(2)} &\in \mathbb{R}^{N \times c} \end{aligned} \right\} \rightarrow X_n = [\mathbf{x}_{n,:}^{(1)}, \mathbf{x}_{n,:}^{(2)}] \in \mathbb{R}^{1 \times (d+c)} \\
& \left. \begin{aligned} Z^{(com)} &\in \mathbb{R}^{N \times K} \\ Z^{(1)} &\in \mathbb{R}^{N \times K_1} \\ Z^{(2)} &\in \mathbb{R}^{N \times K_2} \end{aligned} \right\} \rightarrow \mathbf{z}_{n,:} = [Z_n^{(com)}, Z_n^{(1)}, Z_n^{(2)}] \in \mathbb{R}^{1 \times K_c} \\
& \left. \begin{aligned} A^{(1)} &\in \mathbb{R}^{d \times K} \\ A^{(2)} &\in \mathbb{R}^{c \times K} \\ B^{(1)} &\in \mathbb{R}^{d \times K_1} \\ B^{(2)} &\in \mathbb{R}^{c \times K_2} \end{aligned} \right\} \rightarrow W = \begin{bmatrix} A^{(1)} & B^{(1)} & 0 \\ A^{(2)} & 0 & B^{(2)} \end{bmatrix} \in \mathbb{R}^{(d+c) \times K_c} \\
& \left. \begin{aligned} \Sigma^{(1)} &\in \mathbb{R}^{d \times d} \\ \Sigma^{(2)} &\in \mathbb{R}^{c \times c} \end{aligned} \right\} \rightarrow \Sigma = \begin{bmatrix} \Sigma^{(1)} & 0 \\ 0 & \Sigma^{(2)} \end{bmatrix} \in \mathbb{R}^{(d+c) \times (d+c)}
\end{aligned} \right\} \rightarrow X_n = \mathbf{z}_{n,:} W^T + \epsilon$$
(71)

Let's assume that

$$\mathbf{z}_{n,:} \sim \mathcal{N}(0, I) \quad (72)$$

then, we have that

$$X_n | \mathbf{z}_{n,:} \sim \mathcal{N}(\mathbf{z}_{n,:} W^T, \Sigma) \quad (73)$$

$$\mathbf{x}_{n,:}^{(m)} | \mathbf{z}_{n,:} \sim \mathcal{N}(\mathbf{z}_{n,:} \mathbf{W}^{(m)T}, \Sigma^{(m)}) \quad (74)$$

where  $\mathbf{W}^{(m)}$  is a random variable defined as:

$$\mathbf{W}^{(m)} \sim \text{ARD}(\alpha_0, \beta_0) \quad (75)$$

If we force all columns to have the same prior while letting the rows have a different one this distribution can be written as:

$$\mathbf{w}_{:,k}^{(m)} \sim \mathcal{N}\left(0, \left(\alpha_k^{(m)}\right)^{-1} I\right) \quad (76)$$

This way, high values of  $\alpha_k$  make the column tend to 0 (low variance). On the other hand, if  $\alpha_k$  has a low value you allow the column to have any value (high variance).

$$\alpha_k^{(m)} \sim \text{Gamma}(\alpha_0, \beta_0) \quad (77)$$

where  $\alpha_0$  and  $\beta_0$  are non-critical hyperparameters. This way we are capable of cancelling the components of  $W$  that we don't need. For this reason, the matrix  $W$  is learned disorganised. Intuitively, we have that:

- A 1 in the upper half of the matrix and a 0 in the lower  $\rightarrow B^{(1)}$
- A 0 in the upper half of the matrix and a 1 in the lower  $\rightarrow B^{(2)}$
- A 0 in the upper half of the matrix and a 0 in the lower  $\rightarrow \text{Noise}$
- A 1 in the upper half of the matrix and a 1 in the lower  $\rightarrow A^{(1)}/A^{(2)}$

This can be achieved by forcing sparsity for the latent factors (K) when using ARD (Automatic Relevance Determination).

Conversely, the noise  $\Sigma$  can be modelled as:

$$\Sigma^{(m)} = \tau_m^{-1} I \quad (78)$$

where  $\tau_m$  is,  $\tau_m = \text{Gamma}(\alpha_0^\tau, \beta_0^\tau)$ , being  $\tau$  a model hyperparameter.

### 3.2 Variational inference

Making use of the variational inference definition previously stated, the problem can be defined to apply those techniques to design the model. First of all, considering Figure 3 we can extract the probability distribution of the model

$$p(\Theta | \mathbf{X}^{\{\mathcal{M}\}}) \approx \prod_{m=1}^M \left( q(\mathbf{W}^{(m)}) q(\tau^{(m)}) \prod_{k=1}^{K_c} q(\alpha_k^{(m)}) \right) \prod_{n=1}^N q(\mathbf{z}_{n,:}) \quad (79)$$

This way, we can determine the optimum distribution of the different variables, applying Equation (8) and then calculate the lower bound defined in Equation (6).

In this section the different distributions specified in Equation (79) will be developed making use of the variational inference.

#### 3.2.1 Distribution of $\mathbf{W}^{(m)}$

$$\begin{aligned} \ln(q^*(\mathbf{W}^{(m)})) &= \mathbb{E}_{\mathbf{Z}, \alpha^{(m)}, \tau^{(m)}} \left[ \ln(p(\mathbf{X}^{(m)}, \mathbf{W}^{(m)}, \mathbf{Z}, \alpha^{(m)}, \tau^{(m)})) \right] \\ &= \mathbb{E}_{\mathbf{Z}, \tau^{(m)}} \left[ \ln(p(\mathbf{X}^{(m)} | \mathbf{W}^{(m)}, \mathbf{Z}, \tau^{(m)})) \right] \\ &\quad + \mathbb{E}_{\alpha^{(m)}} \left[ \ln(p(\mathbf{W}^{(m)} | \alpha^{(m)})) \right] + \text{const} \end{aligned} \quad (80)$$

We can now evaluate both terms independently and then sum the results:

$$\begin{aligned} p(\Theta, \tilde{\mathbf{T}}^{\{\mathcal{M}_t\}}, \tilde{\mathbf{X}}^{\{\mathcal{M}_r\}} | \mathbf{T}^{\{\mathcal{M}_t\}}, \mathbf{X}^{\{\mathcal{M}_r\}}) &\approx \\ &= \sum_{n=1}^N \ln \left( \mathcal{N}(\mathbf{z}_{n,:} | \mathbf{W}^{(m)\top}, (\tau^{(m)})^{-1} I) \right) + \text{const} \\ &= \sum_{n=1}^N \left( -\frac{1}{2} \ln \left| (\tau^{(m)})^{-1} I \right| \right. \\ &\quad \left. - \frac{1}{2} (\mathbf{x}_{n,:}^{(m)} - \mathbf{z}_{n,:} | \mathbf{W}^{(m)\top}) \tau^{(m)} (\mathbf{x}_{n,:}^{(m)} - \mathbf{z}_{n,:} | \mathbf{W}^{(m)\top})^T \right) + \text{const} \\ &= -\frac{\tau^{(m)}}{2} \sum_{n=1}^N \left( \mathbf{x}_{n,:}^{(m)} \mathbf{x}_{n,:}^{(m)\top} \right. \\ &\quad \left. - 2 \mathbf{x}_{n,:}^{(m)} \mathbf{W}^{(m)} \mathbf{z}_{n,:}^T + \mathbf{z}_{n,:} \mathbf{W}^{(m)\top} \mathbf{W}^{(m)} \mathbf{z}_{n,:}^T \right) + \text{const} \\ &= \sum_{n=1}^N \left( \tau^{(m)} \mathbf{x}_{n,:}^{(m)} \mathbf{W}^{(m)} \mathbf{z}_{n,:}^T - \frac{\tau^{(m)}}{2} \mathbf{z}_{n,:} \mathbf{W}^{(m)\top} \mathbf{W}^{(m)} \mathbf{z}_{n,:}^T \right) + \text{const} \\ &= \sum_{n=1}^N \left( \tau^{(m)} \mathbf{x}_{n,:}^{(m)} \mathbf{W}^{(m)} \mathbf{z}_{n,:}^T - \frac{\tau^{(m)}}{2} \sum_{d=1}^{D_m} (\mathbf{z}_{n,:} \mathbf{w}_{d,:}^{(m)\top} \mathbf{w}_{d,:}^{(m)} \mathbf{z}_{n,:}^T) \right) + \text{const} \\ &= \sum_{n=1}^N \left( \tau^{(m)} \mathbf{x}_{n,:}^{(m)} \mathbf{W}^{(m)} \mathbf{z}_{n,:}^T - \frac{\tau^{(m)}}{2} \sum_{d=1}^{D_m} (\mathbf{w}_{d,:}^{(m)} \mathbf{z}_{n,:} \mathbf{z}_{n,:} \mathbf{w}_{d,:}^{(m)\top}) \right) + \text{const} \end{aligned} \quad (81)$$

This way the expectation can be calculated as:

$$\mathbb{E}_{\mathbf{Z}, \tau^{(m)}} \left[ \ln \left( p \left( \mathbf{X}^{(m)} \mid \mathbf{W}^{(m)}, \mathbf{Z}, \tau^{(m)} \right) \right) \right] = \sum_{n=1}^N \left( \langle \tau^{(m)} \rangle \mathbf{x}_{n,:}^{(m)} \mathbf{W}^{(m)} \langle \mathbf{z}_{n,:} \rangle - \frac{\langle \tau^{(m)} \rangle}{2} \sum_{d=1}^{D_m} \left( \mathbf{w}_{d,:}^{(m)} \langle \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \rangle \mathbf{w}_{d,:}^{(m)T} \right) \right) \quad (82)$$

Equivalently, the second term would be:

$$\begin{aligned} \ln \left( p \left( \mathbf{W}^{(m)} \mid \boldsymbol{\alpha}^{(m)} \right) \right) &= \sum_{d=1}^{D_m} \sum_{k=1}^{K_c} \ln \left( p \left( \mathbf{w}_{d,k}^{(m)} \mid \alpha_k^{(m)} \right) \right) \\ &= \sum_{d=1}^{D_m} \sum_{k=1}^{K_c} \left( \frac{1}{2} \ln \left( \alpha_k^{(m)} \right) - \frac{\alpha_k^{(m)}}{2} \mathbf{w}_{d,k}^{(m)2} \right) + \text{const} \end{aligned} \quad (83)$$

$$\mathbb{E}_{\boldsymbol{\alpha}^{(m)}} \left[ \ln \left( p \left( \mathbf{W}^{(m)} \mid \boldsymbol{\alpha}^{(m)} \right) \right) \right] = \sum_{d=1}^{D_m} \sum_{k=1}^{K_c} \left( -\frac{\langle \alpha_k^{(m)} \rangle}{2} \mathbf{w}_{d,k}^{(m)2} \right) + \text{const} \quad (84)$$

Finally, joining Equations (82) and (84), we would have that the optimum solution for the variable  $\mathbf{W}^{(m)}$  is:

$$\begin{aligned} \ln \left( q^* \left( \mathbf{W}^{(m)} \right) \right) &= \sum_{d=1}^{D_m} \sum_{k=1}^{K_c} \left( -\frac{\langle \alpha_k^{(m)} \rangle}{2} \mathbf{w}_{d,k}^{(m)2} \right) + \sum_{n=1}^N \left( \langle \tau^{(m)} \rangle \mathbf{x}_{n,:}^{(m)} \mathbf{W}^{(m)} \langle \mathbf{z}_{n,:} \rangle - \frac{\langle \tau^{(m)} \rangle}{2} \sum_{d=1}^{D_m} \left( \mathbf{w}_{d,:}^{(m)} \langle \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \rangle \mathbf{w}_{d,:}^{(m)T} \right) \right) + \text{const} \\ &= \sum_{d=1}^{D_m} \left( -\frac{1}{2} \mathbf{w}_{d,:}^{(m)} \text{diag}(\langle \alpha_k^{(m)} \rangle) \mathbf{w}_{d,:}^{(m)T} + \sum_{n=1}^N \left( \langle \tau^{(m)} \rangle \mathbf{x}_{n,d}^{(m)} \mathbf{w}_{d,:}^{(m)} \langle \mathbf{z}_{n,:} \rangle - \frac{\langle \tau^{(m)} \rangle}{2} \mathbf{w}_{d,:}^{(m)} \langle \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \rangle \mathbf{w}_{d,:}^{(m)T} \right) \right) + \text{const} \\ &= \sum_{d=1}^{D_m} \left( -\frac{1}{2} \mathbf{w}_{d,:}^{(m)} \left( \text{diag}(\langle \alpha_k^{(m)} \rangle) + \langle \tau^{(m)} \rangle \langle \mathbf{z}_{n,:}^T \mathbf{z}_{n,:} \rangle \right) \mathbf{w}_{d,:}^{(m)T} + \langle \tau^{(m)} \rangle \mathbf{w}_{d,:}^{(m)} \langle \mathbf{Z} \rangle^T \mathbf{x}_{:,d}^{(m)T} \right) + \text{const} \end{aligned} \quad (85)$$

This way, comparing the results with the normal distribution, we can identify terms, extracting the following conclusions:

$$q^* \left( \mathbf{W}^{(m)} \right) = \prod_{d=1}^{D_m} \left( \mathcal{N} \left( \mathbf{w}_{d,:}^{(m)} \mid \mu_{\mathbf{w}_{d,:}^{(m)}}, \Sigma_{\mathbf{W}^{(m)}} \right) \right) \quad (86)$$

Where the variance would be:

$$\Sigma_{\mathbf{W}^{(m)}}^{-1} = \text{diag}(\langle \alpha^{(m)} \rangle) + \langle \tau^{(m)} \rangle \langle \mathbf{Z}^T \mathbf{Z} \rangle \quad (87)$$

And the mean could be expressed as:

$$\begin{aligned} \mu_{\mathbf{w}_{d,:}^{(m)}} &= \langle \tau^{(m)} \rangle \Sigma_{\mathbf{W}^{(m)}} \langle \mathbf{Z} \rangle^T \mathbf{x}_{:,d}^{(m)T} \\ \langle \mathbf{W}^{(m)} \rangle &= \langle \tau^{(m)} \rangle \mathbf{X}^{(m)T} \langle \mathbf{Z} \rangle \Sigma_{\mathbf{W}^{(m)}} \end{aligned} \quad (88)$$



### 3.2.2 Distribution of $\mathbf{Z}$

$$\begin{aligned}
\ln(q^*(\mathbf{Z})) &= \mathbb{E}_{\mathbf{W}^{(m)}, \tau^{(m)}} \left[ \ln \left( p \left( \mathbf{X}^{(m)}, \mathbf{W}^{(m)}, \mathbf{Z}, \tau^{(m)} \right) \right) \right] \\
&= \mathbb{E}_{\mathbf{W}^{(1)}, \tau^{(1)}} \left[ \ln \left( p \left( \mathbf{X}^{(1)} | \mathbf{W}^{(1)}, \mathbf{Z}, \tau^{(1)} \right) \right) \right] + \mathbb{E}_{\mathbf{W}^{(2)}, \tau^{(2)}} \left[ \ln \left( p \left( \mathbf{X}^{(2)} | \mathbf{W}^{(2)}, \mathbf{Z}, \tau^{(2)} \right) \right) \right] \\
&\quad + \mathbb{E}[\ln(p(\mathbf{Z}))]
\end{aligned} \tag{89}$$

The first two terms are similar to equation (82), changing the expectation dependency:

$$\begin{aligned}
\mathbb{E}_{\mathbf{W}^{(m)}, \tau^{(m)}} \left[ \ln \left( p \left( \mathbf{X}^{(m)} | \mathbf{W}^{(m)}, \mathbf{Z}, \tau^{(m)} \right) \right) \right] &= \sum_{n=1}^N \left( \langle \tau^{(m)} \rangle \mathbf{x}_{n,:}^{(m)} \langle \mathbf{W}^{(m)} \rangle \mathbf{z}_{n,:}^T \right. \\
&\quad \left. - \frac{\langle \tau^{(m)} \rangle}{2} \mathbf{z}_{n,:} \langle \mathbf{W}^{(m)\top} \mathbf{W}^{(m)} \rangle \mathbf{z}_{n,:}^T \right) + \text{const}
\end{aligned} \tag{90}$$

The third term could be calculated as:

$$\begin{aligned}
\ln(p(\mathbf{Z})) &= \sum_{n=1}^N (\mathcal{N}(0, I)) = \sum_{n=1}^N \left( -\frac{1}{2} \mathbf{z}_{n,:} \mathbf{z}_{n,:}^T \right) + \text{const} \\
\mathbb{E}[\ln(p(\mathbf{Z}))] &= \ln(p(\mathbf{Z})) = \sum_{n=1}^N \left( -\frac{1}{2} \mathbf{z}_{n,:} \mathbf{z}_{n,:}^T \right) + \text{const}
\end{aligned} \tag{91}$$

Therefore, joining the different terms in Equations (90) and (91), the optimum solution has the form:

$$\begin{aligned}
\ln(q^*(\mathbf{Z})) &= \sum_{n=1}^N \left( -\frac{1}{2} \mathbf{z}_{n,:} \mathbf{z}_{n,:}^T + \sum_{m=1}^M \left( \langle \tau^{(m)} \rangle \mathbf{x}_{n,:}^{(m)} \langle \mathbf{W}^{(m)} \rangle \mathbf{z}_{n,:}^T \right. \right. \\
&\quad \left. \left. - \frac{\langle \tau^{(m)} \rangle}{2} \mathbf{z}_{n,:} \langle \mathbf{W}^{(m)\top} \mathbf{W}^{(m)} \rangle \mathbf{z}_{n,:}^T \right) \right) + \text{const}
\end{aligned} \tag{92}$$

This way, we can now identify terms, having that:

$$q^*(\mathbf{Z}) = \prod_{n=1}^N (\mathcal{N}(\mathbf{z}_{n,:} | \mu_{\mathbf{z}_{n,:}}, \Sigma_{\mathbf{Z}})) \tag{93}$$

Where the variance would be:

$$\Sigma_{\mathbf{Z}}^{-1} = I + \sum_{m=1}^M \left( \langle \tau^{(m)} \rangle \langle \mathbf{W}^{(m)\top} \mathbf{W}^{(m)} \rangle \right) \tag{94}$$

And the mean could be expressed as:

$$\begin{aligned}
\mu_{\mathbf{z}_{n,:}} &= \sum_{m=1}^M \left( \langle \tau^{(m)} \rangle \mathbf{x}_{n,:}^{(m)} \langle \mathbf{W}^{(m)} \rangle \Sigma_{\mathbf{Z}} \right) \\
\langle \mathbf{Z} \rangle &= \sum_{m=1}^M \left( \langle \tau^{(m)} \rangle \mathbf{X}^{(m)} \langle \mathbf{W}^{(m)} \rangle \Sigma_{\mathbf{Z}} \right)
\end{aligned} \tag{95}$$

### 3.2.3 Distribution of $\alpha^{(m)}$

$$\ln(q^*(\alpha^{(m)})) = \mathbb{E}_{\mathbf{W}^{(m)}}[\ln(p(\mathbf{W}^{(m)} | \alpha^{(m)}))] + \mathbb{E}[\ln(p(\alpha^{(m)}))] \quad (96)$$

Similarly to what we have done before, we can develop both terms independently:

$$\begin{aligned} \ln(p(\mathbf{W}^{(m)} | \alpha^{(m)})) &= \sum_{d=1}^{D_m} \sum_{k=1}^{K_c} \left( \frac{1}{2} \ln |\alpha_k^{(m)}| - \frac{1}{2} \alpha_k^{(m)} (w_{d,k}^{(m)})^2 \right) + \text{const} \\ &= \sum_{k=1}^{K_c} \left( \frac{D_m}{2} \ln(\alpha_k^{(m)}) - \frac{1}{2} \alpha_k^{(m)} \mathbf{w}_{:,k}^{(m)\top} \mathbf{w}_{:,k}^{(m)} \right) + \text{const} \end{aligned}$$

$$\mathbb{E}[\ln(p(\mathbf{W}^{(m)} | \alpha^{(m)}))] = \sum_{k=1}^{K_c} \left( \frac{D_m}{2} \ln(\alpha_k^{(m)}) - \frac{1}{2} \alpha_k^{(m)} \langle \mathbf{w}_{:,k}^{(m)\top} \mathbf{w}_{:,k}^{(m)} \rangle \right) + \text{const} \quad (97)$$

The second term is deployed as follows:

$$\begin{aligned} \ln(p(\alpha_k^{(m)})) &= -\beta_0 \alpha_k^{(m)} + (\alpha_0 - 1) \ln(\alpha_k^{(m)}) + \text{const} \\ \mathbb{E}[\ln(p(\alpha^{(m)}))] &= \sum_{k=1}^{K_c} (\ln(p(\alpha_k^{(m)}))) = \sum_{k=1}^{K_c} (-\beta_0 \alpha_k^{(m)} + (\alpha_0 - 1) \ln(\alpha_k^{(m)})) + \text{const} \end{aligned} \quad (98)$$

Therefore, by joining both Equations (97) and (98) together, we can now determine the value of the expectation we were looking for

$$\begin{aligned} \ln(q^*(\alpha^{(m)})) &= \sum_{k=1}^{K_c} \left( \frac{D_m}{2} \ln(\alpha_k^{(m)}) - \frac{1}{2} \alpha_k^{(m)} \langle \mathbf{w}_{:,k}^{(m)\top} \mathbf{w}_{:,k}^{(m)} \rangle \right. \\ &\quad \left. - \beta_0 \alpha_k^{(m)} + (\alpha_0 - 1) \ln(\alpha_k^{(m)}) \right) + \text{const} \end{aligned} \quad (99)$$

This way, we can now identify terms, having that:

$$q^*(\alpha^{(m)}) = \prod_{k=1}^{K_c} \left( \text{Gamma}(\alpha_k^{(m)} | a_{\alpha_k^{(m)}}, b_{\alpha_k^{(m)}}) \right) \quad (100)$$

where the variance would be:

$$a_{\alpha_k^{(m)}} = \frac{D_m}{2} + \alpha_0 \quad (101)$$

and the mean could be expressed as:

$$b_{\alpha_k^{(m)}} = \beta_0 + \frac{1}{2} \langle \mathbf{w}_{:,k}^{(m)\top} \mathbf{w}_{:,k}^{(m)} \rangle = \beta_0 + \frac{1}{2} \langle \mathbf{W}^{(m)\top} \mathbf{W}^{(m)} \rangle_{k,k} \quad (102)$$

where

$$\langle \mathbf{w}_{:,k}^{(m)\top} \mathbf{w}_{:,k}^{(m)} \rangle = \sum_{d=1}^{D_m} (w_{d,k}^{(m)})^2 = \langle \mathbf{W}^{(m)\top} \mathbf{W}^{(m)} \rangle_{k,k}$$

### 3.2.4 Distribution of $\tau^{(m)}$

$$\ln(q^*(\tau^{(m)})) = \mathbb{E}_{\mathbf{W}^{(m)}, \mathbf{Z}} [\ln(p(\mathbf{X}^{(m)} | \mathbf{W}^{(m)}, \mathbf{Z}, \tau^{(m)}))] + \mathbb{E} [\ln(p(\tau^{(m)}))] \quad (103)$$

Similarly to what was done in equation (82):

$$\begin{aligned} \ln(p(\mathbf{X}^{(m)} | \mathbf{W}^{(m)}, \mathbf{Z}, \tau^{(m)})) &= \sum_{n=1}^N \ln \left( \mathcal{N} \left( \mathbf{z}_{n,:} \mathbf{W}^{(m)\top}, (\tau^{(m)})^{-1} I \right) \right) + \text{const} \\ &= \sum_{n=1}^N \sum_{d=1}^{D_m} \left( \frac{1}{2} \ln |\tau^{(m)}| \right. \\ &\quad \left. - \frac{\tau^{(m)}}{2} \left( \mathbf{x}_{n,d}^{(m)} - \mathbf{z}_{n,:} \mathbf{w}_{d,:}^{(m)\top} \right)^2 \right) + \text{const} \\ &= \frac{D_m N}{2} \ln(\tau^{(m)}) - \frac{\tau^{(m)}}{2} \sum_{n=1}^N \sum_{d=1}^{D_m} \left( \mathbf{x}_{n,d}^{(m)2} \right. \\ &\quad \left. - 2 \mathbf{w}_{d,:}^{(m)} \mathbf{z}_{n,:}^T \mathbf{x}_{n,d}^{(m)} + \left( \mathbf{z}_{n,:} \mathbf{w}_{d,:}^{(m)\top} \right)^2 \right) + \text{const} \\ &= \frac{D_m N}{2} \ln(\tau^{(m)}) - \frac{\tau^{(m)}}{2} \left( \sum_{n=1}^N \sum_{d=1}^{D_m} \mathbf{x}_{n,d}^{(m)2} \right. \\ &\quad \left. - 2 \sum_{d=1}^{D_m} \mathbf{w}_{d,:}^{(m)} \mathbf{Z}^T \mathbf{x}_{:,d}^{(m)} + \sum_{d=1}^{D_m} \mathbf{w}_{d,:}^{(m)} \mathbf{Z}^T \mathbf{Z} \mathbf{w}_{d,:}^{(m)\top} \right) + \text{const} \\ &= \frac{D_m N}{2} \ln(\tau^{(m)}) - \frac{\tau^{(m)}}{2} \left( \sum_{n=1}^N \sum_{d=1}^{D_m} \mathbf{x}_{n,d}^{(m)2} \right. \\ &\quad \left. - 2 \text{Tr} \{ \mathbf{W}^{(m)} \mathbf{Z}^T \mathbf{X}^{(m)} \} + \text{Tr} \{ \mathbf{W}^{(m)} \mathbf{Z}^T \mathbf{Z} \mathbf{W}^{(m)\top} \} \right) + \text{const} \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{\mathbf{W}^{(m)}, \mathbf{Z}} [\ln(p(\mathbf{X}^{(m)} | \mathbf{W}^{(m)}, \mathbf{Z}, \tau^{(m)}))] &= \frac{D_m N}{2} \ln(\tau^{(m)}) - \frac{\tau^{(m)}}{2} \left( \sum_{n=1}^N \sum_{d=1}^{D_m} \mathbf{x}_{n,d}^{(m)2} \right. \\ &\quad \left. - 2 \text{Tr} \{ \langle \mathbf{W}^{(m)} \rangle \langle \mathbf{Z}^T \rangle \mathbf{X}^{(m)} \} \right. \\ &\quad \left. + \text{Tr} \{ \langle \mathbf{W}^{(m)\top} \mathbf{W}^{(m)} \rangle \langle \mathbf{Z}^T \mathbf{Z} \rangle \} \right) + \text{const} \end{aligned}$$

The second term could subsequently be determined as:

$$\mathbb{E} [\ln(p(\tau^{(m)}))] = \ln(p(\tau^{(m)})) = -\beta_0^\tau \tau^{(m)} + (\alpha_0^\tau - 1) \ln(\tau^{(m)}) + \text{const}$$

Joining them together, we have that:

$$\begin{aligned} \ln(q^*(\tau^{(m)})) &= \frac{D_m N}{2} \ln(\tau^{(m)}) - \frac{\tau^{(m)}}{2} \left( \sum_{n=1}^N \sum_{d=1}^{D_m} \mathbf{x}_{n,d}^{(m)2} \right. \\ &\quad \left. - 2 \text{Tr} \{ \langle \mathbf{W}^{(m)} \rangle \langle \mathbf{Z}^T \rangle \mathbf{X}^{(m)} \} + \text{Tr} \{ \langle \mathbf{W}^{(m)\top} \mathbf{W}^{(m)} \rangle \langle \mathbf{Z}^T \mathbf{Z} \rangle \} \right) \\ &\quad - \beta_0^\tau \tau^{(m)} + (\alpha_0^\tau - 1) \ln(\tau^{(m)}) + \text{const} \end{aligned}$$

So by identifying terms with the distribution, we would have that for this parameter:

$$q^*(\tau^{(m)}) = \left( \text{Gamma}(\tau^{(m)} | a_{\tau^{(m)}}, b_{\tau^{(m)}}) \right) \quad (104)$$

the variance would be:

$$a_{\tau^{(m)}} = \frac{D_m N}{2} + \alpha_0^\tau \quad (105)$$

and the mean could be expressed as:

$$b_{\tau^{(m)}} = \beta_0^\tau + \frac{1}{2} \left( \sum_{n=1}^N \sum_{d=1}^{D_m} x_{n,d}^{(m)2} - 2 \text{Tr} \left\{ \langle \mathbf{W}^{(m)} \rangle \langle \mathbf{Z}^T \rangle \mathbf{X}^{(m)} \right\} + \text{Tr} \left\{ \langle \mathbf{W}^{(m)T} \mathbf{W}^{(m)} \rangle \langle \mathbf{Z}^T \mathbf{Z} \rangle \right\} \right) \quad (106)$$

### 3.3 Update of the lower bound ( $L(q)$ )

Once the optimisation of the different parameters has been determined, we can calculate the changes on the lower bound as:

$$\begin{aligned} L_q &= - \int q(\Theta) \ln \left( \frac{q(\Theta)}{p(X, \Theta)} \right) d\Theta = \int q(\Theta) \ln(p(X, \Theta)) d\Theta - \int q(\Theta) \ln(q(\Theta)) d\Theta \\ &= \mathbb{E}_q[\ln(p(X, \Theta))] - \mathbb{E}_q[\ln(q(\Theta))] \end{aligned} \quad (107)$$

To facilitate the analysis of the different parts of these equations we will evaluate in the following subsections the terms related to  $\mathbb{E}_q[\ln(p(X, \Theta))]$  and to the entropy.

### 3.4 Terms associated to $\mathbb{E}_q[\ln(p(X, \Theta))]$

This first term of the lower bound would be composed by the following terms:

$$\begin{aligned} \mathbb{E}_q[\ln(p(X, \Theta))] &= \mathbb{E}_q[\ln(p(\mathbf{Z}))] + \sum_{m=1}^M \left( \mathbb{E}_q \left[ \ln \left( p(\mathbf{W}^{(m)} | \boldsymbol{\alpha}^{(m)}) \right) \right] + \mathbb{E}_q \left[ \ln \left( p(\boldsymbol{\alpha}^{(m)}) \right) \right] \right) \\ &\quad + \mathbb{E}_q \left[ \ln \left( p(\mathbf{X}^{(m)} | \mathbf{W}^{(m)}, \mathbf{Z}, \tau^{(m)}) \right) \right] + \mathbb{E}_q \left[ \ln \left( p(\tau^{(m)}) \right) \right] \end{aligned} \quad (108)$$

Which, as we have done before, can be independently analysed. To facilitate the extraction of results, these new terms will be included in different subsubsections.

$$\mathbb{E}_q[\ln(p(\mathbf{Z}))]$$

$$\mathbb{E}_q[\ln(p(\mathbf{Z}))] = \sum_{n=1}^N \mathbb{E}_q \left[ -\frac{K_c}{2} \ln(2\pi) - \frac{1}{2} \mathbf{z}_{n,:} \mathbf{z}_{n,:}^T \right] = -\frac{N K_c}{2} \ln(2\pi) - \frac{1}{2} \text{Tr} \{ \langle \mathbf{Z} \mathbf{Z}^T \rangle \} \quad (109)$$

$$\mathbb{E}_q[\ln(p(\mathbf{W}^{(m)} | \boldsymbol{\alpha}^{(m)}))]$$

$$\begin{aligned} \mathbb{E}_q[\ln(p(\mathbf{W}^{(m)} | \boldsymbol{\alpha}^{(m)}))] &= \sum_{k=1}^{K_c} \mathbb{E}_q \left[ -\frac{D_m}{2} \ln(2\pi) + \frac{D_m}{2} \ln(\alpha_k^{(m)}) - \frac{1}{2} \mathbf{w}_{:,k}^{(m)T} \alpha_k^{(m)} \mathbf{w}_{:,k}^{(m)} \right] \\ &= -\frac{K_c D_m}{2} \ln(2\pi) + \frac{D_m}{2} \sum_{k=1}^{K_c} \mathbb{E}_q[\ln(\alpha_k^{(m)})] \\ &\quad - \frac{1}{2} \sum_{k=1}^{K_c} \left( \mathbb{E}_q[\alpha_k^{(m)}] \mathbb{E}_q[\mathbf{w}_{:,k}^{(m)T} \mathbf{w}_{:,k}^{(m)}] \right) \end{aligned} \quad (110)$$

We can now calculate the expectation of the different variables independently to determine the value of the equation, having that:

$$\mathbb{E}_{q(\alpha)}[\ln(\alpha_k^{(m)})] = \psi(a_{\alpha_k^{(m)}}) - \ln(b_{\alpha_k^{(m)}}) \quad (111)$$

$$\mathbb{E}_{q(\alpha)}[\alpha_k^{(m)}] = \frac{a_{\alpha_k^{(m)}}}{b_{\alpha_k^{(m)}}} \quad (112)$$

$$\mathbb{E}_{q(W)}[\mathbf{w}_{:,k}^{(m)T} \mathbf{w}_{:,k}^{(m)}] = \langle \mathbf{w}_{:,k}^{(m)T}, \mathbf{w}_{:,k}^{(m)} \rangle = \langle \mathbf{W}^{(m)T}, \mathbf{W}^{(m)} \rangle_{k,k} \quad (113)$$

We can equate now this to Equation (102) and substituting all these equations on the previous one we have that:

$$\begin{aligned} \mathbb{E}_q[\ln(p(\mathbf{W}^{(m)} | \boldsymbol{\alpha}^{(m)}))] &= -\frac{K_c D_m}{2} \ln(2\pi) + \frac{D_m}{2} \sum_{k=1}^{K_c} \left( \psi(a_{\alpha_k^{(m)}}) - \ln(b_{\alpha_k^{(m)}}) \right) \\ &\quad - \frac{1}{2} \sum_{k=1}^{K_c} \left( \frac{a_{\alpha_k^{(m)}}}{b_{\alpha_k^{(m)}}} (2(b_{\alpha_k^{(m)}} - \beta_0)) \right) \\ &= -\frac{K_c D_m}{2} \ln(2\pi) + \frac{D_m}{2} \sum_{k=1}^{K_c} \left( \psi(a_{\alpha_k^{(m)}}) - \ln(b_{\alpha_k^{(m)}}) \right) \\ &\quad - \sum_{k=1}^{K_c} (a_{\alpha_k^{(m)}}) + \beta_0 \sum_{k=1}^{K_c} \left( \frac{a_{\alpha_k^{(m)}}}{b_{\alpha_k^{(m)}}} \right) \end{aligned} \quad (114)$$

$$\mathbb{E}_q[\ln(p(\boldsymbol{\alpha}^{(m)}))]$$

$$\begin{aligned} \mathbb{E}_q[\ln(p(\boldsymbol{\alpha}^{(m)}))] &= \sum_{k=1}^{K_c} \mathbb{E}_q \left[ \ln(\beta_0) - \beta_0 \alpha_k^{(m)} + (\alpha_0 - 1) \ln(\beta_0) + (\alpha_0 - 1) \ln(\alpha_k^{(m)}) - \ln(\Gamma(\alpha_0)) \right] \\ &= \sum_{k=1}^{K_c} \left( -\beta_0 \mathbb{E}_q[\alpha_k^{(m)}] + \alpha_0 \ln(\beta_0) + (\alpha_0 - 1) \mathbb{E}_q[\ln(\alpha_k^{(m)})] - \ln(\Gamma(\alpha_0)) \right) \end{aligned}$$

Using the expectations determined in Equations (111) and (112) we have:

$$\begin{aligned} \mathbb{E}_q[\ln(p(\boldsymbol{\alpha}^{(m)}))] &= K_c(\alpha_0 \ln(\beta_0) - \ln(\Gamma(\alpha_0))) \\ &\quad + \sum_{k=1}^{K_c} \left( -\beta_0 \frac{a_{\alpha_k^{(m)}}}{b_{\alpha_k^{(m)}}} + (\alpha_0 - 1) \left( \psi(a_{\alpha_k^{(m)}}) - \ln(b_{\alpha_k^{(m)}}) \right) \right) \end{aligned} \quad (115)$$

$$\mathbb{E}_q[\ln(p(\mathbf{W}^{(m)}, \boldsymbol{\alpha}^{(m)}))]$$

We can now calculate the joint distribution of the two previous variables by combining Equations (115) and (114):

$$\begin{aligned} \mathbb{E}_q[\ln(p(\mathbf{W}^{(m)}, \boldsymbol{\alpha}^{(m)}))] &= \left(\frac{D_m}{2} + \alpha_0 - 1\right) \sum_{k=1}^{K_c} \left(\psi(a_{\alpha_k^{(m)}}) - \ln(b_{\alpha_k^{(m)}})\right) \\ &\quad - \frac{K_c D_m}{2} \ln(2\pi) + K_c(\alpha_0 \ln(\beta_0) - \ln(\Gamma(\alpha_0))) - \sum_{k=1}^{K_c} (a_{\alpha_k^{(m)}}) \end{aligned} \quad (116)$$

$$\mathbb{E}_q[\ln(p(\mathbf{X}^{(m)} | \mathbf{W}^{(m)}, \mathbf{Z}, \tau^{(m)}))]$$

$$\begin{aligned} \mathbb{E}_q[\ln(p(\mathbf{X}^{(m)} | \mathbf{W}^{(m)}, \mathbf{Z}, \tau^{(m)}))] &= \sum_{n=1}^N \mathbb{E}_q \left[ -\frac{D_m}{2} \ln(2\pi) + \frac{D_m}{2} \ln(\tau^{(m)}) \right. \\ &\quad \left. - \frac{1}{2} (\mathbf{x}_{n,:}^{(m)} - \mathbf{Z} \mathbf{W}^{(m)T})^T \tau^{(m)} (\mathbf{x}_{n,:}^{(m)} - \mathbf{Z} \mathbf{W}^{(m)T}) \right] \\ &= \sum_{n=1}^N \left( -\frac{D_m}{2} \ln(2\pi) + \frac{D_m}{2} \mathbb{E}_q[\ln(\tau^{(m)})] \right. \\ &\quad \left. - \frac{1}{2} \mathbb{E}_q[\tau^{(m)}] \mathbb{E}_q[\mathbf{x}_{n,:}^{(m)} \mathbf{x}_{n,:}^{(m)T}] \right. \\ &\quad \left. - 2 \mathbf{x}_{n,:}^{(m)} \mathbf{W}^{(m)} \mathbf{z}_{n,:}^T + \mathbf{z}_{n,:} \mathbf{W}^{(m)T} \mathbf{W}^{(m)} \mathbf{z}_{n,:}^T \right] \\ &= -\frac{ND_m}{2} \ln(2\pi) + \frac{D_m}{2} \sum_{n=1}^N (\mathbb{E}_q[\ln(\tau^{(m)})]) \\ &\quad - \frac{1}{2} \mathbb{E}_q[\tau^{(m)}] \sum_{n=1}^N (\mathbf{x}_{n,:}^{(m)} \mathbf{x}_{n,:}^{(m)T}) \\ &\quad - 2 \text{Tr} \left\{ \mathbf{x}_{n,:}^{(m)} \langle \mathbf{W}^{(m)} \rangle \langle \mathbf{z}_{n,:}^T \rangle \right\} + \text{Tr} \left\{ \langle \mathbf{W}^{(m)T}, \mathbf{W}^{(m)} \rangle \langle \mathbf{z}_{n,:}^T, \mathbf{z}_{n,:} \rangle \right\} \end{aligned} \quad (117)$$

If we compare and substitute the value obtained in Equation (106):

$$\begin{aligned} \mathbb{E}_q[\ln(p(\mathbf{X}^{(m)} | \mathbf{W}^{(m)}, \mathbf{Z}, \tau^{(m)}))] &= -\frac{ND_m}{2} \ln(2\pi) + \frac{D_m}{2} (\psi(a_{\tau^{(m)}}) - \ln(b_{\tau^{(m)}})) - \frac{a_{\tau^{(m)}}}{b_{\tau^{(m)}}} (b_{\tau^{(m)}} - \beta_0^\tau) \\ &= -\frac{ND_m}{2} \ln(2\pi) + \frac{D_m}{2} (\psi(a_{\tau^{(m)}}) - \ln(b_{\tau^{(m)}})) - a_{\tau^{(m)}} + \frac{a_{\tau^{(m)}}}{b_{\tau^{(m)}}} \beta_0^\tau \end{aligned} \quad (118)$$

$$\mathbb{E}_q[\ln(p(\tau^{(m)}))]$$

$$\begin{aligned} \mathbb{E}_q[\ln(p(\tau^{(m)}))] &= \mathbb{E}_q[\ln(\beta_0^\tau) - \beta_0^\tau \tau^{(m)} + (\alpha_0^\tau - 1) \ln(\beta_0^\tau) + (\alpha_0^\tau - 1) \ln(\tau^{(m)}) - \ln(\Gamma(\alpha_0^\tau))] \\ &= -\beta_0^\tau \mathbb{E}_q[\tau^{(m)}] + \alpha_0^\tau \ln(\beta_0^\tau) + (\alpha_0^\tau - 1) \mathbb{E}_q[\ln(\tau^{(m)})] - \ln(\Gamma(\alpha_0^\tau)) \\ &= \alpha_0^\tau \ln(\beta_0^\tau) - \ln(\Gamma(\alpha_0^\tau)) - \beta_0^\tau \frac{a_{\tau^{(m)}}}{b_{\tau^{(m)}}} + (\alpha_0^\tau - 1)(\psi(a_{\tau^{(m)}}) - \ln(b_{\tau^{(m)}})) \end{aligned} \quad (119)$$

$$\mathbb{E}_q[\ln(p(\mathbf{X}^{(m)}, \tau^{(m)} | \mathbf{W}^{(m)}, \mathbf{Z}))]$$

In a similar way to what was done before, we can now compute the joint distribution of  $\mathbf{X}^{(m)}$  and  $\tau^{(m)}$  by combining Equations (119) and (118) we have:

$$\begin{aligned} \mathbb{E}_q[\ln(p(\mathbf{X}^{(m)}, \tau^{(m)} | \mathbf{W}^{(m)}, \mathbf{Z}))] &= -\frac{ND_m}{2} \ln(2\pi) - a_{\tau^{(m)}} + \alpha_0^\tau \ln(\beta_0^\tau) - \ln(\Gamma(\alpha_0^\tau)) + \\ &\quad \left( \frac{D_m}{2} + \alpha_0^\tau - 1 \right) (\psi(a_{\tau^{(m)}}) - \ln(b_{\tau^{(m)}})) \end{aligned} \quad (120)$$

### 3.5 Terms of entropy, $\mathbb{E}_q[\ln(q(\Theta))]$

This first term of the lower bound would be composed by the following terms:

$$\begin{aligned} \mathbb{E}_q[\ln(q(\Theta))] &= \mathbb{E}_q[\ln(q(\mathbf{Z}))] + \sum_{m=1}^M \left( \mathbb{E}_q[\ln(q(\mathbf{W}^{(m)}))] \right. \\ &\quad \left. + \mathbb{E}_q[\ln(q(\boldsymbol{\alpha}^{(m)}))] + \mathbb{E}_q[\ln(q(\tau^{(m)}))] \right) \end{aligned}$$

Which, as we have done before, can be independently analysed. To facilitate the extraction of results, these new terms will be included in different subsections.

$$\mathbb{E}_q[\ln(q(\mathbf{Z}))]$$

$$\begin{aligned} \mathbb{E}_q[\ln(q(\mathbf{Z}))] &= \sum_{n=1}^N \left( \frac{K_c}{2} \ln(2\pi e) + \frac{1}{2} \ln |\Sigma_{\mathbf{Z}}| \right) \\ &= \frac{NK_c}{2} \ln(2\pi e) + \frac{N}{2} \ln |\Sigma_{\mathbf{Z}}| \end{aligned} \quad (121)$$

$$\mathbb{E}_q[\ln(q(\mathbf{W}^{(m)}))]$$

$$\begin{aligned} \mathbb{E}_q[\ln(q(\mathbf{W}^{(m)}))] &= \sum_{d=1}^{D_m} \left( \frac{K_c}{2} \ln(2\pi e) + \frac{1}{2} \ln |\Sigma_{\mathbf{W}^{(m)}}| \right) \\ &= \frac{D_m K_c}{2} \ln(2\pi e) + \frac{D_m}{2} \ln |\Sigma_{\mathbf{W}^{(m)}}| \end{aligned} \quad (122)$$

$$\mathbb{E}_q[\ln(q(\boldsymbol{\alpha}^{(m)}))]$$

$$\mathbb{E}_q[\ln(q(\boldsymbol{\alpha}^{(m)}))] = \sum_{k=1}^{K_c} \left( a_{\alpha_k^{(m)}} + \ln(\Gamma(a_{\alpha_k^{(m)}})) - (1 - a_{\alpha_k^{(m)}}) \psi(a_{\alpha_k^{(m)}}) - \ln(b_{\alpha_k^{(m)}}) \right) \quad (123)$$

$$\mathbb{E}_q[\ln(q(\tau^{(m)}))]$$

$$\mathbb{E}_q[\ln(q(\tau^{(m)}))] = a_{\tau^{(m)}} + \ln(\Gamma(a_{\tau^{(m)}})) - (1 - a_{\tau^{(m)}}) \psi(a_{\tau^{(m)}}) - \ln(b_{\tau^{(m)}}) \quad (124)$$

### 3.6 Complete definition of the lower bound ( $L(q)$ )

Once the different terms have been calculated, we can now determine how the lower bound is going to update by joining together all the previously calculated terms:

$$\begin{aligned} L_q &= \mathbb{E}_q[\ln(p(X, \Theta))] - \mathbb{E}_q[\ln(q(\Theta))] \\ &= \mathbb{E}_q[\ln(p(\mathbf{Z}))] + \sum_{m=1}^M \left( \mathbb{E}_q[\ln(p(\mathbf{W}^{(m)}, \boldsymbol{\alpha}^{(m)}))] \right. \\ &\quad \left. + \mathbb{E}_q[\ln(p(\mathbf{X}^{(m)}, \tau^{(m)} | \mathbf{W}^{(m)}, \mathbf{Z}))] \right) \\ &\quad - \mathbb{E}_q[\ln(q(\mathbf{Z}))] - \sum_{m=1}^M \left( \mathbb{E}_q[\ln(q(\mathbf{W}^{(m)}))] + \mathbb{E}_q[\ln(q(\boldsymbol{\alpha}^{(m)}))] \right. \\ &\quad \left. + \mathbb{E}_q[\ln(q(\tau^{(m)}))] \right) \end{aligned}$$



$$\begin{aligned}
= & -\frac{NK_c}{2} \ln(2\pi) - \frac{1}{2} \text{Tr}\{\langle \mathbf{Z}^T \mathbf{Z} \rangle\} \\
& + \sum_{m=1}^M \left( -\frac{K_c D_m}{2} \ln(2\pi) + \left( \frac{D_m}{2} + \alpha_0 - 1 \right) \sum_{k=1}^{K_c} \left( \psi(a_{\alpha_k^{(m)}}) - \ln(b_{\alpha_k^{(m)}}) \right) \right. \\
& \left. + K_c (\alpha_0 \ln(\beta_0) - \ln(\Gamma(\alpha_0))) - \sum_{k=1}^{K_c} \left( a_{\alpha_k^{(m)}} \right) \right) \\
& + \sum_{m=1}^M \left( -\frac{ND_m}{2} \ln(2\pi) - a_{\tau^{(m)}} + \alpha_0^\tau \ln(\beta_0^\tau) - \ln(\Gamma(\alpha_0^\tau)) \right. \\
& \left. + \left( \frac{D_m}{2} + \alpha_0^\tau - 1 \right) (\psi(a_{\tau^{(m)}}) - \ln(b_{\tau^{(m)}})) \right) \\
& - \frac{NK_c}{2} \ln(2\pi e) - \frac{N}{2} \ln |\Sigma_{\mathbf{Z}}| \\
& - \sum_{m=1}^M \left( \frac{D_m K_c}{2} \ln(2\pi e) + \frac{D_m}{2} \ln |\Sigma_{\mathbf{W}^{(m)}}| \right) \\
& - \sum_{m=1}^M \left( \sum_{k=1}^{K_c} \left( a_{\alpha_k^{(m)}} + \ln(\Gamma(a_{\alpha_k^{(m)}})) \right) - (1 - a_{\alpha_k^{(m)}}) \psi(a_{\alpha_k^{(m)}}) - \ln(b_{\alpha_k^{(m)}}) \right) \\
& - \sum_{m=1}^M \left( a_{\tau^{(m)}} + \ln(\Gamma(a_{\tau^{(m)}})) - (1 - a_{\tau^{(m)}}) - \ln(b_{\tau^{(m)}}) \right) \tag{125}
\end{aligned}$$

As this equation is going to be iterated through to check the state of the bound, we can simplify it by eliminating all the terms that are constant through the iteration process:

$$\begin{aligned}
L_q = & -\frac{1}{2} \text{Tr}\{\langle \mathbf{Z}^T \mathbf{Z} \rangle\} \\
& - \sum_{m=1}^M \left( \left( \frac{D_m}{2} + \alpha_0 - 1 \right) \sum_{k=1}^{K_c} (\ln(b_{\alpha_k^{(m)}})) \right) \\
& - \sum_{m=1}^M \left( \left( \frac{D_m}{2} + \alpha_0^\tau - 1 \right) (\ln(b_{\tau^{(m)}})) \right) \\
& - \frac{N}{2} \ln |\Sigma_{\mathbf{Z}}| - \sum_{m=1}^M \left( \frac{D_m}{2} \ln |\Sigma_{\mathbf{W}^{(m)}}| \right) + \sum_{m=1}^M \left( \sum_{k=1}^{K_c} (\ln(b_{\alpha_k^{(m)}})) \right) \\
& + \sum_{m=1}^M (\ln(b_{\tau^{(m)}})) \tag{126}
\end{aligned}$$

## References

Bishop, C.M., 1999. Bayesian pca. *Advances in neural information processing systems*, 382–388.