

Probabilistic latent variable models for shape correspondence analysis

Hernán F. García and Mauricio A. Álvarez

October 6, 2017

Abstract

This report provides a detailed analysis of the mathematical aspects related to our model for shape correspondence analysis.

1 The model

Based on the Iwatta’s paper (see [1, 2]), we develop a non linear model for shape correspondence analysis using probabilistic latent variable models.

Suppose that we are given objects in D domains $\mathcal{X} = \{\mathbf{X}_d\}_{d=1}^D$ mapped to a Hilbert space \mathcal{H} , where $\mathbf{X}_d = \{\mathbf{x}_{dn}\}_{n=1}^{N_d}$ is a set of objects in the d th domain, and $\mathbf{x}_{dn} \in \mathbb{R}^{M_d}$ is the feature vector of the n th object in the d th domain. By introducing a function $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ called the *kernel*, that performs a given mapping over the objects, $\phi : \mathcal{X} \rightarrow \mathcal{H}$ such that $\forall x, x' \in \mathcal{X}$,

$$k(x, x') := \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}, \quad (1)$$

we can cluster groups of correspondences by using a non-linear function that represents the shape descriptors in the Hilbert space.

Our notation is summarized in Table 1. Note that we are unaware of any correspondence between objects in different domains. The number of objects N_d and the dimensionality L_d for each domain in \mathcal{H} can be different from those of other domains. Therefore, our task is to match clusters of descriptors (groupwise correspondences) across multiple brain structures in an unsupervised manner [2].

As in infinite Gaussian mixture models, our approach assumes that there are an infinite number of clusters related to each correspondence, and each cluster j has a latent vector $\boldsymbol{\zeta}_j \in \mathbb{R}^P$ in a latent space of dimension P . Descriptors that have the same cluster assignments s_{dn} are related by the same latent vector and considered to match (establish a groupwise correspondence).

Each object in $\phi(\mathbf{x}_{dn}) \in \mathcal{H}$ in the d th domain is generated depending on the domain-specific projection matrix $\mathbf{W}_d = [\varphi(\mathbf{w}_1^d), \varphi(\mathbf{w}_2^d), \dots, \varphi(\mathbf{w}_P^d)]$, $\mathbf{W}_d \in \mathbb{R}^{L_d \times P}$ and latent vector $\boldsymbol{\zeta}_{s_{dn}}$ that is selected from a set of latent vectors $\mathbf{Z} = \{\boldsymbol{\zeta}_j\}_{j=1}^{\infty}$. Here, $s_{dn} = \{1, \dots, \infty\}$ is the latent cluster assignment of object $\phi(\mathbf{x}_{dn})$.

The proposed model is based on an infinite mixture model, where the probability of descriptor mapped in a Hilbert space $\phi(\mathbf{x}_{dn})$ is given by

Table 1: Notation.

Symbol in \mathcal{I}	Symbol in \mathcal{H}	Description
D		Number of shapes
N_d		Number of objects (3D-landmarks) in the d th shape
M_d	L_d	Dimensionality of the observed features in the d th shape
K	P	Dimensionality of a latent vector
J	Q	Number of correspondences (latent vectors) to which objects are assigned
\mathbf{x}_{dn}	$\phi(\mathbf{x}_{dn})$	Observation of the n th object in the d th shape, $\mathbf{x}_{dn} \in \mathbb{R}^{M_d}$
\mathbf{z}_j	$\boldsymbol{\zeta}_j$	Latent vector for the j th correspondence, $\boldsymbol{\zeta}_j \in \mathbb{R}^K$
\mathbf{W}_d	\mathbf{W}_d	Projection matrix for the d th domain, $\mathbf{W}_d \in \mathbb{R}^{M_d \times K}$
θ_j		Mixture weight for the j th cluster, $\theta_j \geq 0$, $\sum_{j=1}^{\infty} \theta_j = 1$

$$p(\phi(\mathbf{x}_{dn})|\mathbf{Z}, \mathbf{W}, \boldsymbol{\theta}) = \sum_{j=1}^{\infty} \theta_j \mathcal{N}(\phi(\mathbf{x}_{dn}) | \mathbf{W}_d \boldsymbol{\zeta}_j, \alpha^{-1} \mathbf{I}), \quad (2)$$

where $\mathbf{W} = \{\mathbf{W}_d\}_{d=1}^D$ is a set of projections matrices, $\boldsymbol{\theta} = (\theta_j)_{j=1}^{\infty}$ are the mixture weights, θ_j represents the probability that the j th cluster is chosen, and $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes a normal distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. One important contribution derived in [1], is that we can analyze multiples structures with different properties and dimensionalities, by employing projection matrices for each brain structure (domain-specific). Figure 1 shows the scheme of the proposed model, in which the relationship between shape descriptors, and latent vectors is described.

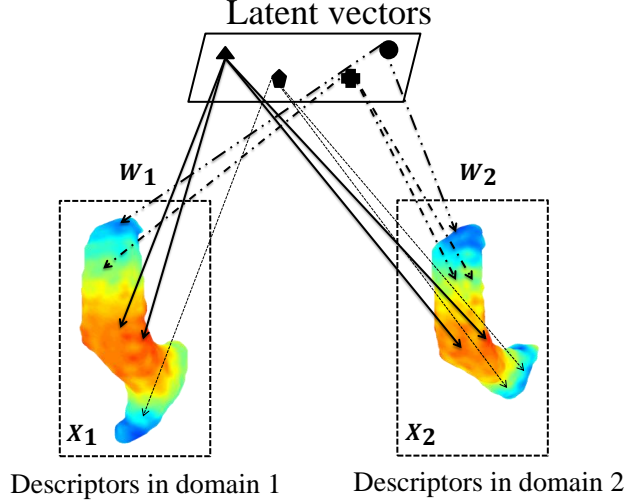


Figure 1: Scheme for the groupwise correspondence method. The figure shows an example of establishing clusters of correspondences in two domains (left ventrals).

In order to draw the cluster proportions, we use a stick-breaking process to generate mixture weights for a Dirichlet process with concentration parameter γ [1] (a, b and r are the hyperparameters). The joint probability of the data \mathcal{X} , and the cluster assignments $\mathbf{S} = \left\{ \{s_{dn}\}_{n=1}^{N_d} \right\}_{d=1}^D$ are given by

$$p(\boldsymbol{\Phi}, \mathbf{S} | \mathbf{W}, a, b, r, \gamma) = p(\mathbf{S} | \gamma) p(\boldsymbol{\Phi} | \mathbf{S}, \mathbf{W}, a, b, r). \quad (3)$$

By marginalizing out the mixture weights $\boldsymbol{\theta}$, $p(\mathbf{S} | \gamma)$ becomes in

$$p(\mathbf{S} | \gamma) = \frac{\gamma^J \prod_{j=1}^J (N_{.j} - 1)!}{\gamma (\gamma + 1) \dots (\gamma + N - 1)}, \quad (4)$$

where $N = \sum_{d=1}^D N_d$ is the total number of shape descriptors, $N_{.j}$ represents the number of descriptors assigned to the cluster j , and J is the number of clusters that satisfies $N_{.j} > 0$.

1.1 The linear model (likelihood)

By marginalizing out latent vectors \mathbf{Z} and the precision parameter α , the second factor of (??) is computed by

$$\begin{aligned}
p(\mathbf{X}|\mathbf{S}, \mathbf{W}, a, b, r) &= \int \int \prod_{d=1}^D \prod_{n=1}^{N_d} \mathcal{N}(\mathbf{x}_{dn} | \mathbf{W}_d \mathbf{z}_{s_{dn}}, \alpha^{-1} \mathbf{I}) \mathcal{G}(\alpha | a, b) \times \prod_{j=1}^J \mathcal{N}(\mathbf{z}_j | \mathbf{0}, (\alpha r)^{-1} \mathbf{I}) d\mathbf{Z} d\alpha \\
&= \int \int \prod_{d=1}^D \prod_{n=1}^{N_d} \left(\frac{\alpha}{2\pi} \right)^{M_d/2} \exp\left(-\frac{\alpha}{2} \|\mathbf{x}_{dn} - \mathbf{W}_d \mathbf{z}_{s_{dn}}\|^2\right) \prod_{j=1}^J \left(\frac{\alpha r}{2\pi} \right)^{K/2} \\
&\quad \times \exp\left(-\frac{\alpha r}{2} \|\mathbf{z}_j\|^2\right) \frac{b^a \alpha^{a-1}}{\Gamma(a)} \exp(-b\alpha) d\mathbf{Z} d\alpha \\
&= \frac{b^a}{\Gamma(a)} \int \int \left(\frac{\alpha}{2\pi} \right)^{\sum_d M_d N_d/2} \exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \|\mathbf{x}_{dn} - \mathbf{W}_d \mathbf{z}_{s_{dn}}\|^2\right) \left(\frac{\alpha r}{2\pi} \right)^{KJ/2} \\
&\quad \times \exp\left(-\frac{\alpha r}{2} \sum_{j=1}^J \|\mathbf{z}_j\|^2\right) \exp(-b\alpha) \alpha^{a-1} d\mathbf{Z} d\alpha
\end{aligned} \tag{5}$$

Solving for the first exponential term

$$\begin{aligned}
\exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \|\mathbf{x}_{dn} - \mathbf{W}_d \mathbf{z}_{s_{dn}}\|^2\right) &= \exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} [\mathbf{x}_{dn}^\top \mathbf{x}_{dn} \right. \\
&\quad \left. - 2\mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{x}_{dn} + \mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{W}_d \mathbf{z}_{s_{dn}}]\right)
\end{aligned} \tag{6}$$

The equation in (5) becomes

$$\begin{aligned}
p(\mathbf{X}|\mathbf{S}, \mathbf{W}, a, b, r) &= \frac{b^a}{\Gamma(a)} \int \int \left(\frac{\alpha}{2\pi} \right)^{\sum_d M_d N_d/2} \exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} [\mathbf{x}_{dn}^\top \mathbf{x}_{dn} \right. \\
&\quad \left. - 2\mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{x}_{dn} + \mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{W}_d \mathbf{z}_{s_{dn}}]\right) \left(\frac{\alpha r}{2\pi} \right)^{KJ/2} \\
&\quad \times \exp\left(-\frac{\alpha r}{2} \sum_{j=1}^J \mathbf{z}_j^\top \mathbf{z}_j\right) \exp(-b\alpha) \alpha^{a-1} d\mathbf{Z} d\alpha.
\end{aligned} \tag{7}$$

The exponential terms in (7) becomes

$$\begin{aligned}
&\exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} [\mathbf{x}_{dn}^\top \mathbf{x}_{dn} - 2\mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{x}_{dn} + \mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{W}_d \mathbf{z}_{s_{dn}}] - \frac{\alpha r}{2} \sum_{j=1}^J \mathbf{z}_j^\top \mathbf{z}_j - b\alpha\right) = \\
&\exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} [\mathbf{x}_{dn}^\top \mathbf{x}_{dn} - b\alpha]\right) \exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} [-2\mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{x}_{dn} + \mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{W}_d \mathbf{z}_{s_{dn}}] \right. \\
&\quad \left. - \frac{\alpha r}{2} \sum_{j=1}^J \mathbf{z}_j^\top \mathbf{z}_j\right)
\end{aligned} \tag{8}$$

By analyzing the n th objects that has the cluster assignment j ($n : s_{dn} = j$), the second factor in (8) becomes

$$\begin{aligned}
& \exp \left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \left[-2\mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{x}_{dn} + \mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{W}_d \mathbf{z}_{s_{dn}} \right] - \frac{\alpha r}{2} \sum_{j=1}^J \mathbf{z}_j^\top \mathbf{z}_j \right) = \exp \left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n:s_{dn} \neq j} \left[-2\mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{x}_{dn} \right. \right. \\
& \left. \left. + \mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{W}_d \mathbf{z}_{s_{dn}} \right] \right) \times \underbrace{\exp \left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n:s_{dn}=j} \left[-2\mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{x}_{dn} + \mathbf{z}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{W}_d \mathbf{z}_{s_{dn}} \right] - \frac{\alpha r}{2} \sum_{j=1}^J \mathbf{z}_j^\top \mathbf{z}_j \right)}_C \\
& = \exp \left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{j=1}^J N_{dj} \left[-2\mathbf{z}_j^\top \mathbf{W}_d^\top \mathbf{x}_{dn} + \mathbf{z}_j^\top \mathbf{W}_d^\top \mathbf{W}_d \mathbf{z}_j \right] - \frac{\alpha r}{2} \sum_{j=1}^J \mathbf{z}_j^\top \mathbf{z}_j \right) \\
& = \exp \left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{j=1}^J \left[-2\mathbf{z}_j^\top \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} + \mathbf{z}_j^\top N_{dj} \mathbf{W}_d^\top \mathbf{W}_d \mathbf{z}_j \right] - \frac{\alpha r}{2} \sum_{j=1}^J \mathbf{z}_j^\top \mathbf{z}_j \right) \\
& = \exp \left(-\frac{\alpha}{2} \sum_{j=1}^J \left[-2\mathbf{z}_j^\top \sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} + \mathbf{z}_j^\top \sum_{d=1}^D N_{dj} \mathbf{W}_d^\top \mathbf{W}_d \mathbf{z}_j \right] - \frac{\alpha r}{2} \sum_{j=1}^J \mathbf{z}_j^\top \mathbf{z}_j \right) \\
& = \exp \left(-\frac{\alpha}{2} \sum_{j=1}^J \left[-2\mathbf{z}_j^\top \sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} + \mathbf{z}_j^\top \sum_{d=1}^D N_{dj} \mathbf{W}_d^\top \mathbf{W}_d \mathbf{z}_j + r \mathbf{z}_j^\top \mathbf{z}_j \right] \right) \\
& = \exp \left(-\frac{\alpha}{2} \sum_{j=1}^J \left[-2\mathbf{z}_j^\top \sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} + \mathbf{z}_j^\top \left(\sum_{d=1}^D N_{dj} \mathbf{W}_d^\top \mathbf{W}_d + r \mathbf{I} \right) \mathbf{z}_j \right] \right). \tag{9}
\end{aligned}$$

By using the quadratic property

$$-\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^\top \mathbf{C}^{-1} (\mathbf{z} - \boldsymbol{\mu}) = -\frac{1}{2} [\mathbf{z}^\top \mathbf{C}^{-1} \mathbf{z} - 2\mathbf{z}^\top \mathbf{C}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^\top \mathbf{C}^{-1} \boldsymbol{\mu}], \tag{10}$$

where

$$\mathbf{C}_j^{-1} = \sum_{d=1}^D N_{dj} \mathbf{W}_d^\top \mathbf{W}_d + r \mathbf{I}, \tag{11}$$

and

$$\begin{aligned}
-2\mathbf{z}^\top \mathbf{C}^{-1} \boldsymbol{\mu} &= -2\mathbf{z}_j^\top \sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \\
\boldsymbol{\mu}_j &= \mathbf{C}_j \sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn}.
\end{aligned} \tag{12}$$

By completing the square as: $\arg = \arg + \frac{1}{2} \boldsymbol{\mu}^\top \mathbf{C}^{-1} \boldsymbol{\mu} - \frac{1}{2} \boldsymbol{\mu}^\top \mathbf{C}^{-1} \boldsymbol{\mu}$, the argument in (9) becomes

$$\begin{aligned}
& \exp \left(-\frac{\alpha}{2} \sum_{j=1}^J \left[-2\mathbf{z}_j^\top \sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} + \mathbf{z}_j^\top \left(\sum_{d=1}^D N_{dj} \mathbf{W}_d^\top \mathbf{W}_d + r \mathbf{I} \right) \mathbf{z}_j \right] \right) \\
& = \exp \left(-\frac{\alpha}{2} \left[\sum_{j=1}^J (\mathbf{z}_j - \boldsymbol{\mu}_j)^\top \mathbf{C}_j^{-1} (\mathbf{z}_j - \boldsymbol{\mu}_j) \right] \right) \exp \left(-\frac{\alpha}{2} \sum_{j=1}^J \boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \boldsymbol{\mu}_j \right)
\end{aligned} \tag{13}$$

Substituting (13) in (7) give us

$$\begin{aligned}
p(\mathbf{X}|\mathbf{S}, \mathbf{W}, a, b, r) &= \frac{b^a}{\Gamma(a)} \iint \left(\frac{\alpha}{2\pi} \right)^{\sum_d M_d N_d / 2} \left(\frac{\alpha r}{2\pi} \right)^{KJ/2} \exp \left(-\frac{\alpha}{2} \left[\sum_{j=1}^J (\mathbf{z}_j - \boldsymbol{\mu}_j)^\top \mathbf{C}_j^{-1} (\mathbf{z}_j - \boldsymbol{\mu}_j) \right] \right) d\mathbf{Z} \\
& \exp \left(-\alpha \left[\frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \mathbf{x}_{dn}^\top \mathbf{x}_{dn} - \frac{1}{2} \sum_{j=1}^J \boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \boldsymbol{\mu}_j + b \right] \right) \alpha^{a-1} d\alpha
\end{aligned} \tag{14}$$

In equation (14), factors related to \mathbf{Z} are grouped together. We integrated out \mathbf{Z} using

$$\int \exp \left(-\frac{1}{2} (\mathbf{z}_j - \boldsymbol{\mu}_j)^\top [\alpha^{-1} \mathbf{C}_j]^{-1} (\mathbf{z}_j - \boldsymbol{\mu}_j) \right) d\mathbf{z}_j = (2\pi)^{K/2} |\alpha^{-1} \mathbf{C}_j|^{1/2} = (2\pi)^{K/2} \alpha^{-K/2} |\mathbf{C}_j|^{1/2}, \quad (15)$$

which is the normalization constant of P -dimensional Gaussian distribution. Since we have the sum over the number of correspondences (latent vectors), K , the above equation ranges for all of these clusters. The equation (14), becomes

$$p(\mathbf{X}|\mathbf{S}, \mathbf{W}, a, b, r) = \frac{b^a}{\Gamma(a)} \int \left(\frac{\alpha}{2\pi} \right)^{\sum_d M_d N_d / 2} \left(\frac{\alpha r}{2\pi} \right)^{KJ/2} \prod_{j=1}^J \left[(2\pi)^{K/2} \alpha^{-K/2} |\mathbf{C}_j|^{1/2} \right] \exp \left(-\alpha \left[\frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \mathbf{x}_{dn}^\top \mathbf{x}_{dn} - \frac{1}{2} \sum_{j=1}^J \boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \boldsymbol{\mu}_j + b \right] \right) \alpha^{a-1} d\alpha \quad (16)$$

$$= \frac{b^a}{\Gamma(a)} \int \left(\frac{\alpha}{2\pi} \right)^{\sum_d M_d N_d / 2} \left(\frac{\alpha r}{2\pi} \right)^{KJ/2} (2\pi)^{KJ/2} \alpha^{-KJ/2} \prod_{j=1}^J |\mathbf{C}_j|^{1/2} \exp \left(-\alpha \left[\frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \mathbf{x}_{dn}^\top \mathbf{x}_{dn} - \frac{1}{2} \sum_{j=1}^J \boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \boldsymbol{\mu}_j + b \right] \right) \alpha^{a-1} d\alpha \quad (17)$$

The α parameter is integrated out by using the following normalization constant of a Gamma distribution

$$\int \alpha^{a'-1} \exp(-b\alpha) d\alpha = \frac{\Gamma(a')}{b^{a'}}. \quad (18)$$

Finally the likelihood is given by

$$p(\mathbf{X}|\mathbf{S}, \mathbf{W}, a, b, r) = (2\pi)^{-\frac{\sum_d M_d N_d}{2}} r^{\frac{KJ}{2}} \frac{b^a}{b^{a'}} \frac{\Gamma(a')}{\Gamma(a)} \prod_{j=1}^J |\mathbf{C}_j|^{1/2}, \quad (19)$$

Here,

$$a' = a + \frac{\sum_d M_d N_d}{2}, \quad b' = b + \frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \mathbf{x}_{dn}^\top \mathbf{x}_{dn} - \frac{1}{2} \sum_{j=1}^J \boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \boldsymbol{\mu}_j, \quad (20)$$

$$\boldsymbol{\mu}_j = \mathbf{C}_j \sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn}, \quad \mathbf{C}_j^{-1} = \sum_{d=1}^D N_{dj} \mathbf{W}_d^\top \mathbf{W}_d + r\mathbf{I}, \quad (21)$$

where N_{dj} is the number of descriptors assigned to cluster j in the shape d (domain).

1.2 Inference for the Linear Model

1.2.1 M-Step Linear

In the M-step, the projection matrices \mathbf{W} are estimated by maximizing the logarithm of the joint likelihood (3). The gradient of the joint likelihood is computed by

$$\frac{\partial \log p(\mathbf{X}, \mathbf{S}|\mathbf{W}, a, b, r, \gamma)}{\partial \mathbf{W}_d} = \frac{\partial \log p(\mathbf{S}|\gamma)}{\partial \mathbf{W}_d} + \frac{\partial \log p(\mathbf{X}|\mathbf{S}, \mathbf{W}, a, b, r)}{\partial \mathbf{W}_d}, \quad (22)$$

Since the derivative of the first term in the above expression is zero, the expression becomes

$$\begin{aligned}
\frac{\partial \log p(\mathbf{X}, \mathbf{S} | \mathbf{W}, a, b, r, \gamma)}{\partial \mathbf{W}_d} &= \frac{\partial \log p(\mathbf{X} | \mathbf{S}, \mathbf{W}, a, b, r)}{\partial \mathbf{W}_d} = \frac{\partial \log \left[(2\pi)^{-\frac{\sum_d M_d N_d}{2}} r^{\frac{JK}{2}} \frac{b^a}{b'^{a'}} \frac{\Gamma(a')}{\Gamma(a)} \prod_{j=1}^J |\mathbf{C}_j|^{1/2} \right]}{\partial \mathbf{W}_d}, \\
&= \frac{\partial \log \left[\frac{cte}{b'^{a'}} \prod_{j=1}^J |\mathbf{C}_j|^{1/2} \right]}{\partial \mathbf{W}_d}, \\
&= -\frac{a'}{b'} \frac{\partial b'}{\partial \mathbf{W}_d} + \frac{1}{2} \sum_{j=1}^J \frac{\partial \log |\mathbf{C}_j|}{\partial \mathbf{W}_d} + 0, \\
&= -\frac{a'}{b'} \frac{\partial b'}{\partial \mathbf{W}_d} + \frac{1}{2} \sum_{j=1}^J \text{tr} \left(\mathbf{C}_j^{-1} \frac{\partial \mathbf{C}_j}{\partial \mathbf{W}_d} \right)
\end{aligned} \tag{23}$$

where

$$\frac{\partial b'}{\partial \mathbf{W}_d} = \frac{\partial}{\partial \mathbf{W}_d} \left[\frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \mathbf{x}_{dn}^\top \mathbf{x}_{dn} - \frac{1}{2} \sum_{j=1}^J \boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \boldsymbol{\mu}_j + b \right], \tag{24}$$

Here the second factor of the argument is the only which depends on \mathbf{W}_d .

$$\begin{aligned}
\frac{\partial b'}{\partial \mathbf{W}_d} &= \frac{\partial}{\partial \mathbf{W}_d} \left[-\frac{1}{2} \sum_{j=1}^J \boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \boldsymbol{\mu}_j \right] \\
&= -\frac{1}{2} \sum_{j=1}^J \frac{\partial}{\partial \mathbf{W}_d} [\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \boldsymbol{\mu}_j] = -\frac{1}{2} \sum_{j=1}^J \text{tr} \left(\frac{\partial}{\partial \mathbf{W}_d} [\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \boldsymbol{\mu}_j] \right) \\
&= -\frac{1}{2} \sum_{j=1}^J \text{tr} \left(\frac{\partial \boldsymbol{\mu}_j^\top}{\partial \mathbf{W}_d} \mathbf{C}_j^{-1} \boldsymbol{\mu}_j + \boldsymbol{\mu}_j^\top \left(\frac{\partial \mathbf{C}_j^{-1}}{\partial \mathbf{W}_d} \boldsymbol{\mu}_j + \mathbf{C}_j^{-1} \frac{\partial \boldsymbol{\mu}_j}{\partial \mathbf{W}_d} \right) \right) \\
&= -\frac{1}{2} \sum_{j=1}^J \text{tr} (\partial \boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \boldsymbol{\mu}_j) + \text{tr} (\boldsymbol{\mu}_j^\top \partial \mathbf{C}_j^{-1} \boldsymbol{\mu}_j) + \text{tr} (\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \partial \boldsymbol{\mu}_j),
\end{aligned} \tag{25}$$

by applying trace properties (transpose elements)

$$\begin{aligned}
\frac{\partial b'}{\partial \mathbf{W}_d} &= -\frac{1}{2} \sum_{j=1}^J \text{tr} (\partial \boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \boldsymbol{\mu}_j) + \text{tr} (\boldsymbol{\mu}_j^\top \partial \mathbf{C}_j^{-1} \boldsymbol{\mu}_j) + \text{tr} (\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \partial \boldsymbol{\mu}_j) \\
&= -\frac{1}{2} \sum_{j=1}^J \text{tr} (\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \partial \boldsymbol{\mu}_j) + \text{tr} (\boldsymbol{\mu}_j^\top \partial \mathbf{C}_j^{-1} \boldsymbol{\mu}_j) + \text{tr} (\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \partial \boldsymbol{\mu}_j) \\
&= -\frac{1}{2} \sum_{j=1}^J \text{tr} (\boldsymbol{\mu}_j^\top \partial \mathbf{C}_j^{-1} \boldsymbol{\mu}_j) + 2 \text{tr} (\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \partial \boldsymbol{\mu}_j) \\
&= \underbrace{-\sum_{j=1}^J \text{tr} (\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \partial \boldsymbol{\mu}_j)}_A - \underbrace{\frac{1}{2} \sum_{j=1}^J \text{tr} (\boldsymbol{\mu}_j^\top \partial \mathbf{C}_j^{-1} \boldsymbol{\mu}_j)}_B
\end{aligned} \tag{26}$$

First, for the B part we have:

$$\begin{aligned}
-\frac{1}{2} \sum_{j=1}^J \text{tr}(\boldsymbol{\mu}_j^\top \partial \mathbf{C}_j^{-1} \boldsymbol{\mu}_j) &= -\frac{1}{2} \sum_{j=1}^J \text{tr} \left(\boldsymbol{\mu}_j^\top \partial \left[\sum_{d=1}^D N_{dj} \mathbf{W}_d^\top \mathbf{W}_d + r \mathbf{I} \right] \boldsymbol{\mu}_j \right) \\
&= -\frac{1}{2} \sum_{j=1}^J \text{tr}(\boldsymbol{\mu}_j^\top [N_{dj} (\partial \mathbf{W}_d^\top \mathbf{W}_d + \mathbf{W}_d^\top \partial \mathbf{W}_d)] \boldsymbol{\mu}_j) \\
&= -\frac{1}{2} \sum_{j=1}^J [N_{dj} \text{tr}(\boldsymbol{\mu}_j^\top \partial \mathbf{W}_d^\top \mathbf{W}_d \boldsymbol{\mu}_j) + N_{dj} \text{tr}(\boldsymbol{\mu}_j^\top \mathbf{W}_d^\top \partial \mathbf{W}_d \boldsymbol{\mu}_j)] \\
&= -\frac{1}{2} \sum_{j=1}^J [N_{dj} \text{tr}((\boldsymbol{\mu}_j^\top \partial \mathbf{W}_d^\top \mathbf{W}_d \boldsymbol{\mu}_j)^\top) + N_{dj} \text{tr}(\boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \mathbf{W}_d^\top \partial \mathbf{W}_d)] \\
&= -\frac{1}{2} \sum_{j=1}^J [N_{dj} \text{tr}(\boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \mathbf{W}_d^\top \partial \mathbf{W}_d \boldsymbol{\mu}_j) + N_{dj} \text{tr}(\boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \mathbf{W}_d^\top \partial \mathbf{W}_d)] \\
&= -\sum_{j=1}^J N_{dj} \text{tr}(\boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \mathbf{W}_d^\top \partial \mathbf{W}_d). \tag{27}
\end{aligned}$$

By using derivatives properties for trace forms as

$$\frac{\partial \text{tr}[F(\mathbf{X})]}{\partial \mathbf{X}} = f(\mathbf{X})^\top, \tag{28}$$

where $f(\cdot)$ is the scalar derivative of $F(\cdot)$, the equation (27) becomes

$$-\frac{1}{2} \sum_{j=1}^J \text{tr}(\boldsymbol{\mu}_j^\top \partial \mathbf{C}_j^{-1} \boldsymbol{\mu}_j) = -\sum_{j=1}^J N_{dj} (\boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \mathbf{W}_d^\top)^\top = -\sum_{j=1}^J N_{dj} \mathbf{W}_d \boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top. \tag{29}$$

Besides, for the A part we have:

$$-\sum_{j=1}^J \text{tr}(\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \partial \boldsymbol{\mu}_j) \rightarrow \frac{\partial \boldsymbol{\mu}_j}{\partial \mathbf{W}_d} = \frac{\partial}{\partial \mathbf{W}_d} \left[\mathbf{C}_j \sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \right]. \tag{30}$$

The derivative for $\boldsymbol{\mu}_j$ is given by

$$\frac{\partial \boldsymbol{\mu}_j}{\partial \mathbf{W}_d} = \underbrace{\partial \mathbf{C}_j \left(\sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \right)}_C + \underbrace{\mathbf{C}_j \partial \left(\sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \right)}_D \tag{31}$$

For the C part, we have

$$\begin{aligned}
\partial \mathbf{C}_j \left(\sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \right) &= \frac{\partial}{\partial \mathbf{W}_d} \left[\left(\sum_{d=1}^D N_{dj} \mathbf{W}_d^\top \mathbf{W}_d + r \mathbf{I} \right)^{-1} \right] \left(\sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \right) \\
&= -\mathbf{C}_j \frac{\partial}{\partial \mathbf{W}_d} \left[\sum_{d=1}^D N_{dj} \mathbf{W}_d^\top \mathbf{W}_d + r \mathbf{I} \right] \mathbf{C}_j \sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \\
&= -\mathbf{C}_j (N_{dj} [\partial \mathbf{W}_d^\top \mathbf{W}_d + \mathbf{W}_d^\top \partial \mathbf{W}_d]) \boldsymbol{\mu}_j. \tag{32}
\end{aligned}$$

The D part is computed as

$$\mathbf{C}_j \partial \left(\sum_{d=1}^D \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \right) = \mathbf{C}_j \partial \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn}. \tag{33}$$

Then the A part becomes,

$$\begin{aligned}
-\sum_{j=1}^J \text{tr}(\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \partial \boldsymbol{\mu}_j) &= -\sum_{j=1}^J \text{tr} \left(\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \left[-\mathbf{C}_j (N_{dj} [\partial \mathbf{W}_d^\top \mathbf{W}_d + \mathbf{W}_d^\top \partial \mathbf{W}_d]) \boldsymbol{\mu}_j + \mathbf{C}_j \partial \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \right] \right) \\
&= \sum_{j=1}^J N_{dj} \text{tr}(\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \mathbf{C}_j (\partial \mathbf{W}_d^\top \mathbf{W}_d) \boldsymbol{\mu}_j) + \sum_{j=1}^J N_{dj} \text{tr}(\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \mathbf{C}_j (\mathbf{W}_d^\top \partial \mathbf{W}_d) \boldsymbol{\mu}_j) \\
&\quad - \sum_{j=1}^J \text{tr} \left(\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \mathbf{C}_j \partial \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \right) \\
&= \sum_{j=1}^J N_{dj} \text{tr}(\boldsymbol{\mu}_j^\top \partial \mathbf{W}_d^\top \mathbf{W}_d \boldsymbol{\mu}_j) + \sum_{j=1}^J N_{dj} \text{tr}(\boldsymbol{\mu}_j^\top \mathbf{W}_d^\top \partial \mathbf{W}_d \boldsymbol{\mu}_j) \\
&\quad - \sum_{j=1}^J \text{tr} \left(\boldsymbol{\mu}_j^\top \partial \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \right) \\
&= \sum_{j=1}^J N_{dj} \text{tr} \left((\boldsymbol{\mu}_j^\top \partial \mathbf{W}_d^\top \mathbf{W}_d \boldsymbol{\mu}_j)^\top \right) + \sum_{j=1}^J N_{dj} \text{tr}(\boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \mathbf{W}_d^\top \partial \mathbf{W}_d) \\
&\quad - \sum_{j=1}^J \text{tr} \left(\left(\boldsymbol{\mu}_j^\top \partial \mathbf{W}_d^\top \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \right)^\top \right) \\
&= \sum_{j=1}^J N_{dj} \text{tr}(\boldsymbol{\mu}_j^\top \mathbf{W}_d^\top \partial \mathbf{W}_d \boldsymbol{\mu}_j) + \sum_{j=1}^J N_{dj} \text{tr}(\boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \mathbf{W}_d^\top \partial \mathbf{W}_d) \\
&\quad - \sum_{j=1}^J \text{tr} \left(\sum_{n:s_{dn}=j} \mathbf{x}_{dn}^\top \partial \mathbf{W}_d \boldsymbol{\mu}_j \right) \\
&= \sum_{j=1}^J N_{dj} \text{tr}(\boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \mathbf{W}_d^\top \partial \mathbf{W}_d) + \sum_{j=1}^J N_{dj} \text{tr}(\boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \mathbf{W}_d^\top \partial \mathbf{W}_d) \\
&\quad - \sum_{j=1}^J \text{tr} \left(\boldsymbol{\mu}_j \sum_{n:s_{dn}=j} \mathbf{x}_{dn}^\top \partial \mathbf{W}_d \right) \\
&= 2 \sum_{j=1}^J N_{dj} \text{tr}(\boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \mathbf{W}_d^\top \partial \mathbf{W}_d) - \sum_{j=1}^J \text{tr} \left(\boldsymbol{\mu}_j \sum_{n:s_{dn}=j} \mathbf{x}_{dn}^\top \partial \mathbf{W}_d \right). \tag{34}
\end{aligned}$$

By using the derivatives properties for trace forms described above, we have

$$\begin{aligned}
-\sum_{j=1}^J \text{tr}(\boldsymbol{\mu}_j^\top \mathbf{C}_j^{-1} \partial \boldsymbol{\mu}_j) &= 2 \sum_{j=1}^J N_{dj} (\boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \mathbf{W}_d^\top)^\top - \sum_{j=1}^J \left(\boldsymbol{\mu}_j \sum_{n:s_{dn}=j} \mathbf{x}_{dn}^\top \right)^\top \\
&= 2 \sum_{j=1}^J N_{dj} \mathbf{W}_d \boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top - \sum_{j=1}^J \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \boldsymbol{\mu}_j^\top. \tag{35}
\end{aligned}$$

Finally

$$\begin{aligned}
\frac{\partial b'}{\partial \mathbf{W}_d} &= 2 \sum_{j=1}^J N_{dj} \mathbf{W}_d \boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top - \sum_{j=1}^J \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \boldsymbol{\mu}_j^\top - \sum_{j=1}^J N_{dj} \mathbf{W}_d \boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \\
&= \sum_{j=1}^J N_{dj} \mathbf{W}_d \boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top - \sum_{j=1}^J \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \boldsymbol{\mu}_j^\top \\
&= \sum_{j=1}^J \left\{ N_{dj} \mathbf{W}_d \boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top - \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \boldsymbol{\mu}_j^\top \right\}.
\end{aligned} \tag{36}$$

For the part B ,

$$\begin{aligned}
\frac{1}{2} \sum_{j=1}^J \text{Tr} \left(\mathbf{C}_j^{-1} \frac{\partial \mathbf{C}_j}{\partial \mathbf{W}_d} \right) &= \frac{1}{2} \sum_{j=1}^J \text{tr} \left(\mathbf{C}_j^{-1} \left(-\mathbf{C}_j \left(N_{dj} [\partial \mathbf{W}_d^\top \mathbf{W}_d + \mathbf{W}_d^\top \partial \mathbf{W}_d] \right) \mathbf{C}_j \right) \right) \\
&= -\frac{1}{2} \sum_{j=1}^J N_{dj} \text{tr} \left(\mathbf{C}_j^{-1} \mathbf{C}_j [\partial \mathbf{W}_d^\top \mathbf{W}_d + \mathbf{W}_d^\top \partial \mathbf{W}_d] \mathbf{C}_j \right) \\
&= -\frac{1}{2} \sum_{j=1}^J N_{dj} \text{tr} \left(\mathbf{C}_j^{-1} \mathbf{C}_j [\partial \mathbf{W}_d^\top \mathbf{W}_d + \mathbf{W}_d^\top \partial \mathbf{W}_d] \mathbf{C}_j \right) \\
&= -\frac{1}{2} \sum_{j=1}^J N_{dj} \text{tr} (\partial \mathbf{W}_d^\top \mathbf{W}_d \mathbf{C}_j) - \frac{1}{2} \sum_{j=1}^J N_{dj} \text{tr} (\mathbf{W}_d^\top \partial \mathbf{W}_d \mathbf{C}_j) \\
&= -\frac{1}{2} \sum_{j=1}^J N_{dj} \text{tr} \left((\partial \mathbf{W}_d^\top \mathbf{W}_d \mathbf{C}_j)^\top \right) - \frac{1}{2} \sum_{j=1}^J N_{dj} \text{tr} (\mathbf{C}_j \mathbf{W}_d^\top \partial \mathbf{W}_d) \\
&= -\frac{1}{2} \sum_{j=1}^J N_{dj} \text{tr} (\mathbf{C}_j \mathbf{W}_d^\top \partial \mathbf{W}_d) - \frac{1}{2} \sum_{j=1}^J N_{dj} \text{tr} (\mathbf{C}_j \mathbf{W}_d^\top \partial \mathbf{W}_d) \\
&= -\sum_{j=1}^J N_{dj} \text{tr} (\mathbf{C}_j \mathbf{W}_d^\top \partial \mathbf{W}_d) = -\sum_{j=1}^J N_{dj} \mathbf{W}_d \mathbf{C}_j.
\end{aligned} \tag{37}$$

Finally the derivative of the log-likelihood is computed as

$$\begin{aligned}
\frac{\partial \log p(\mathbf{X}, \mathbf{S} | \mathbf{W}, a, b, r, \gamma)}{\partial \mathbf{W}_d} &= -\frac{a'}{b'} \frac{\partial b'}{\partial \mathbf{W}_d} + \frac{1}{2} \sum_{j=1}^J \text{tr} \left(\mathbf{C}_j^{-1} \frac{\partial \mathbf{C}_j}{\partial \mathbf{W}_d} \right) \\
&= -\frac{a'}{b'} \left[\sum_{j=1}^J \left\{ N_{dj} \mathbf{W}_d \boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top - \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \boldsymbol{\mu}_j^\top \right\} \right] - \sum_{j=1}^J N_{dj} \mathbf{W}_d \mathbf{C}_j.
\end{aligned} \tag{38}$$

We can obtain the projection matrices that maximize the joint likelihood analytically as follows,

$$\mathbf{W}_d = -\frac{a'}{b'} \left(\sum_{j=1}^J \sum_{n:s_{dn}=j} \mathbf{x}_{dn} \boldsymbol{\mu}_j^\top \right) \left(\sum_{j=1}^J N_{dj} \mathbf{C}_j + \frac{a'}{b'} N_{dj} \boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top \right)^{-1}. \tag{39}$$

1.3 Likelihood for the non-linear model

For the non-linear model, we give the derivation of the likelihood in (2), in which latent vectors \mathbf{Z} and precision parameter α are analytically integrated out. To this end, we use two mappings functions $\phi(\mathbf{x}_{dn})$ and $\varphi(\mathbf{w}_j^d)$, in order to represent our observations in the Hilbert space. First let us define the following expressions:

Let us define the projection matrix \mathbf{W}_d in \mathcal{H} as

$$\mathbf{W}_d = \begin{pmatrix} \varphi_1(\mathbf{w}_1^d) & \varphi_1(\mathbf{w}_2^d) & \cdots & \varphi_1(\mathbf{w}_P^d) \\ \varphi_2(\mathbf{w}_1^d) & \varphi_2(\mathbf{w}_2^d) & \cdots & \varphi_2(\mathbf{w}_P^d) \\ \vdots & \vdots & \ddots & \vdots \\ \varphi_{L_d}(\mathbf{w}_1^d) & \varphi_{L_d}(\mathbf{w}_2^d) & \cdots & \varphi_{L_d}(\mathbf{w}_P^d) \end{pmatrix} \quad (40)$$

$$\begin{aligned} p(\Phi|\mathbf{S}, \mathbf{W}, a, b, r) &= \int \int \prod_{d=1}^D \prod_{n=1}^{N_d} \mathcal{N}(\phi(\mathbf{x}_{dn}) | \mathbf{W}_d \boldsymbol{\zeta}_{s_{dn}}, \alpha^{-1} \mathbf{I}) \mathcal{G}(\alpha | a, b) \times \prod_{j=1}^Q \mathcal{N}(\boldsymbol{\zeta}_j | \mathbf{0}, (\alpha r)^{-1} \mathbf{I}) d\mathbf{Z} d\alpha \\ &= \int \int \prod_{d=1}^D \prod_{n=1}^{N_d} \left(\frac{\alpha}{2\pi} \right)^{L_d/2} \exp\left(-\frac{\alpha}{2} \|\phi(\mathbf{x}_{dn}) - \mathbf{W}_d \boldsymbol{\zeta}_{s_{dn}}\|^2\right) \prod_{j=1}^Q \left(\frac{\alpha r}{2\pi} \right)^{P/2} \\ &\quad \times \exp\left(-\frac{\alpha r}{2} \|\boldsymbol{\zeta}_j\|^2\right) \frac{b^a \alpha^{a-1}}{\Gamma(a)} \exp(-b\alpha) d\mathbf{Z} d\alpha \\ &= \frac{b^a}{\Gamma(a)} \int \int \left(\frac{\alpha}{2\pi} \right)^{\sum_d L_d N_d / 2} \exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \|\phi(\mathbf{x}_{dn}) - \mathbf{W}_d \boldsymbol{\zeta}_{s_{dn}}\|^2\right) \left(\frac{\alpha r}{2\pi} \right)^{PQ/2} \\ &\quad \times \exp\left(-\frac{\alpha r}{2} \sum_{j=1}^Q \|\boldsymbol{\zeta}_j\|^2\right) \exp(-b\alpha) \alpha^{a-1} d\mathbf{Z} d\alpha \end{aligned} \quad (41)$$

Solving for the first exponential term

$$\begin{aligned} \exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \|\phi(\mathbf{x}_{dn}) - \mathbf{W}_d \boldsymbol{\zeta}_{s_{dn}}\|^2\right) &= \exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} [\phi(\mathbf{x}_{dn})^\top \phi(\mathbf{x}_{dn}) \right. \\ &\quad \left. - 2\boldsymbol{\zeta}_{s_{dn}}^\top \mathbf{W}_d^\top \phi(\mathbf{x}_{dn}) + \boldsymbol{\zeta}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{W}_d \boldsymbol{\zeta}_{s_{dn}}]\right) \end{aligned} \quad (42)$$

The equation in (41) becomes

$$\begin{aligned} p(\Phi|\mathbf{S}, \mathbf{W}, a, b, r) &= \frac{b^a}{\Gamma(a)} \int \int \left(\frac{\alpha}{2\pi} \right)^{\sum_d L_d N_d / 2} \exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} [\phi(\mathbf{x}_{dn})^\top \phi(\mathbf{x}_{dn}) \right. \\ &\quad \left. - 2\boldsymbol{\zeta}_{s_{dn}}^\top \mathbf{W}_d^\top \phi(\mathbf{x}_{dn}) + \boldsymbol{\zeta}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{W}_d \boldsymbol{\zeta}_{s_{dn}}]\right) \left(\frac{\alpha r}{2\pi} \right)^{PQ/2} \\ &\quad \times \exp\left(-\frac{\alpha r}{2} \sum_{j=1}^Q \boldsymbol{\zeta}_j^\top \boldsymbol{\zeta}_j\right) \exp(-b\alpha) \alpha^{a-1} d\mathbf{Z} d\alpha. \end{aligned} \quad (43)$$

The exponential terms in (7) becomes

$$\begin{aligned} &\exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} [\phi(\mathbf{x}_{dn})^\top \phi(\mathbf{x}_{dn}) - 2\boldsymbol{\zeta}_{s_{dn}}^\top \mathbf{W}_d^\top \phi(\mathbf{x}_{dn}) + \boldsymbol{\zeta}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{W}_d \boldsymbol{\zeta}_{s_{dn}}] - \frac{\alpha r}{2} \sum_{j=1}^Q \boldsymbol{\zeta}_j^\top \boldsymbol{\zeta}_j - b\alpha\right) = \\ &\exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} [\phi(\mathbf{x}_{dn})^\top \phi(\mathbf{x}_{dn}) - b\alpha]\right) \exp\left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} [-2\boldsymbol{\zeta}_{s_{dn}}^\top \mathbf{W}_d^\top \phi(\mathbf{x}_{dn}) + \boldsymbol{\zeta}_{s_{dn}}^\top \mathbf{W}_d^\top \mathbf{W}_d \boldsymbol{\zeta}_{s_{dn}}] \right. \\ &\quad \left. - \frac{\alpha r}{2} \sum_{j=1}^Q \boldsymbol{\zeta}_j^\top \boldsymbol{\zeta}_j\right) \end{aligned} \quad (44)$$

By analyzing the n th objects that has the cluster assignment j ($n : s_{dn} = j$), the second factor in (44) becomes

$$\begin{aligned}
& \exp \left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \left[-2\zeta_{s_{dn}}^\top \mathbf{w}_d^\top \phi(\mathbf{x}_{dn}) + \zeta_{s_{dn}}^\top \mathbf{w}_d^\top \mathbf{w}_d \zeta_{s_{dn}} \right] - \frac{\alpha r}{2} \sum_{j=1}^Q \zeta_j^\top \zeta_j \right) \\
&= \exp \left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n:s_{dn} \neq j} \left[-2\zeta_{s_{dn}}^\top \mathbf{w}_d^\top \phi(\mathbf{x}_{dn}) + \zeta_{s_{dn}}^\top \mathbf{w}_d^\top \mathbf{w}_d \zeta_{s_{dn}} \right] \right) \\
&\times \exp \left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{n:s_{dn}=j} \left[-2\zeta_{s_{dn}}^\top \mathbf{w}_d^\top \phi(\mathbf{x}_{dn}) + \zeta_{s_{dn}}^\top \mathbf{w}_d^\top \mathbf{w}_d \zeta_{s_{dn}} \right] - \frac{\alpha r}{2} \sum_{j=1}^Q \zeta_j^\top \zeta_j \right) \\
&= \exp \left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{j=1}^Q N_{dj} \left[-2\zeta_j^\top \mathbf{w}_d^\top \phi(\mathbf{x}_{dn}) + \zeta_j^\top \mathbf{w}_d^\top \mathbf{w}_d \zeta_j \right] - \frac{\alpha r}{2} \sum_{j=1}^Q \zeta_j^\top \zeta_j \right) \\
&= \exp \left(-\frac{\alpha}{2} \sum_{d=1}^D \sum_{j=1}^Q \left[-2\zeta_j^\top \mathbf{w}_d^\top \sum_{n:s_{dn}=j} \phi(\mathbf{x}_{dn}) + \zeta_j^\top N_{dj} \mathbf{w}_d^\top \mathbf{w}_d \zeta_j \right] - \frac{\alpha r}{2} \sum_{j=1}^Q \zeta_j^\top \zeta_j \right) \\
&= \exp \left(-\frac{\alpha}{2} \sum_{j=1}^Q \left[-2\zeta_j^\top \sum_{d=1}^D \mathbf{w}_d^\top \sum_{n:s_{dn}=j} \phi(\mathbf{x}_{dn}) + \zeta_j^\top \sum_{d=1}^D N_{dj} \mathbf{w}_d^\top \mathbf{w}_d \zeta_j \right] - \frac{\alpha r}{2} \sum_{j=1}^Q \zeta_j^\top \zeta_j \right) \\
&= \exp \left(-\frac{\alpha}{2} \sum_{j=1}^Q \left[-2\zeta_j^\top \sum_{d=1}^D \mathbf{w}_d^\top \sum_{n:s_{dn}=j} \phi(\mathbf{x}_{dn}) + \zeta_j^\top \sum_{d=1}^D N_{dj} \mathbf{w}_d^\top \mathbf{w}_d \zeta_j + r \zeta_j^\top \zeta_j \right] \right) \\
&= \exp \left(-\frac{\alpha}{2} \sum_{j=1}^Q \left[-2\zeta_j^\top \sum_{d=1}^D \mathbf{w}_d^\top \sum_{n:s_{dn}=j} \phi(\mathbf{x}_{dn}) + \zeta_j^\top \left(\sum_{d=1}^D N_{dj} \mathbf{w}_d^\top \mathbf{w}_d + r \mathbf{I} \right) \zeta_j \right] \right), \tag{45}
\end{aligned}$$

where N_{dj} is the number of objects assigned to cluster j in the domain d . By using the quadratic property

$$-\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^\top \boldsymbol{\Lambda}^{-1} (\mathbf{z} - \boldsymbol{\mu}) = -\frac{1}{2} [\mathbf{z}^\top \boldsymbol{\Lambda}^{-1} \mathbf{z} - 2\mathbf{z}^\top \boldsymbol{\Lambda}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^\top \boldsymbol{\Lambda}^{-1} \boldsymbol{\mu}], \tag{46}$$

where

$$\boldsymbol{\Lambda}_j^{-1} = \sum_{d=1}^D N_{dj} \underbrace{\mathbf{w}_d^\top \mathbf{w}_d}_A + r \mathbf{I}, \tag{47}$$

and

$$\begin{aligned}
-2\mathbf{z}^\top \boldsymbol{\Lambda}^{-1} \boldsymbol{\mu} &= -2\zeta_j^\top \sum_{d=1}^D \mathbf{w}_d^\top \sum_{n:s_{dn}=j} \phi(\mathbf{x}_{dn}) \\
\boldsymbol{\mu}_j &= \boldsymbol{\Lambda}_j \sum_{d=1}^D \mathbf{w}_d^\top \sum_{n:s_{dn}=j} \phi(\mathbf{x}_{dn}) \\
\boldsymbol{\mu}_j &= \boldsymbol{\Lambda}_j \sum_{d=1}^D \sum_{n:s_{dn}=j} \underbrace{\mathbf{w}_d^\top \phi(\mathbf{x}_{dn})}_B. \tag{48}
\end{aligned}$$

We can rewrite the expression in A by placing a given kernel for \mathbf{w}_d

$$\begin{aligned}
\mathcal{W}_d^\top \mathcal{W}_d &= \begin{bmatrix} \varphi(\mathbf{w}_1^d)^\top \varphi(\mathbf{w}_1^d) & \varphi(\mathbf{w}_1^d)^\top \varphi(\mathbf{w}_2^d) & \cdots & \varphi(\mathbf{w}_1^d)^\top \varphi(\mathbf{w}_P^d) \\ \varphi(\mathbf{w}_2^d)^\top \varphi(\mathbf{w}_1^d) & \varphi(\mathbf{w}_2^d)^\top \varphi(\mathbf{w}_2^d) & \cdots & \varphi(\mathbf{w}_2^d)^\top \varphi(\mathbf{w}_P^d) \\ \vdots & \vdots & \ddots & \vdots \\ \varphi(\mathbf{w}_P^d)^\top \varphi(\mathbf{w}_1^d) & \varphi(\mathbf{w}_P^d)^\top \varphi(\mathbf{w}_2^d) & \cdots & \varphi(\mathbf{w}_P^d)^\top \varphi(\mathbf{w}_P^d) \end{bmatrix} \\
&= \begin{bmatrix} k(\mathbf{w}_1^d, \mathbf{w}_1^d) & k(\mathbf{w}_1^d, \mathbf{w}_2^d) & \cdots & k(\mathbf{w}_1^d, \mathbf{w}_P^d) \\ k(\mathbf{w}_2^d, \mathbf{w}_1^d) & k(\mathbf{w}_2^d, \mathbf{w}_2^d) & \cdots & k(\mathbf{w}_2^d, \mathbf{w}_P^d) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{w}_P^d, \mathbf{w}_1^d) & k(\mathbf{w}_P^d, \mathbf{w}_2^d) & \cdots & k(\mathbf{w}_P^d, \mathbf{w}_P^d) \end{bmatrix} = \hat{\mathbf{K}}_d.
\end{aligned} \tag{49}$$

We can rewrite the expression in B by placing a given kernel of the form

$$\mathcal{W}_d^\top \phi(\mathbf{x}_{dn}) = \begin{bmatrix} \varphi(\mathbf{w}_1^d)^\top \phi(\mathbf{x}_{dn}) \\ \varphi(\mathbf{w}_2^d)^\top \phi(\mathbf{x}_{dn}) \\ \vdots \\ \varphi(\mathbf{w}_P^d)^\top \phi(\mathbf{x}_{dn}) \end{bmatrix} = \begin{bmatrix} k(\mathbf{w}_1^d, \mathbf{x}_{dn}) \\ k(\mathbf{w}_2^d, \mathbf{x}_{dn}) \\ \vdots \\ k(\mathbf{w}_P^d, \mathbf{x}_{dn}) \end{bmatrix} = \hat{\mathbf{k}}_d, \tag{50}$$

where $\hat{\mathbf{k}}_d$ represents the kernel evaluated in the \mathbf{x}_{dn} objects that has the cluster assignment j .

By completing the square as: $\arg = \arg + \frac{1}{2} \boldsymbol{\mu}^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu} - \frac{1}{2} \boldsymbol{\mu}^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}$, the argument in (9) becomes

$$\begin{aligned}
&\exp \left(-\frac{\alpha}{2} \sum_{j=1}^Q \left[-2 \boldsymbol{\zeta}_j^\top \sum_{d=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d + \boldsymbol{\zeta}_j^\top \left(\sum_{d=1}^D N_{dj} \hat{\mathbf{K}}_d + r \mathbf{I} \right) \boldsymbol{\zeta}_j \right] \right) \\
&= \exp \left(-\frac{\alpha}{2} \left[\sum_{j=1}^Q (\boldsymbol{\zeta}_j - \boldsymbol{\mu}_j)^\top \boldsymbol{\Lambda}_j^{-1} (\boldsymbol{\zeta}_j - \boldsymbol{\mu}_j) \right] \right) \exp \left(-\frac{\alpha}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right)
\end{aligned} \tag{51}$$

Substituting (51) in (43) give us

$$\begin{aligned}
p(\boldsymbol{\Phi} | \mathbf{S}, \mathcal{W}, a, b, r) &= \frac{b^a}{\Gamma(a)} \iint \left(\frac{\alpha}{2\pi} \right)^{\sum_d L_d N_d / 2} \left(\frac{\alpha r}{2\pi} \right)^{PQ/2} \exp \left(-\frac{\alpha}{2} \left[\sum_{j=1}^Q (\boldsymbol{\zeta}_j - \boldsymbol{\mu}_j)^\top \boldsymbol{\Lambda}_j^{-1} (\boldsymbol{\zeta}_j - \boldsymbol{\mu}_j) \right] \right) d\mathbf{Z} \\
&\exp \left(-\alpha \left[\frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) - \frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j + b \right] \right) \alpha^{a-1} d\alpha
\end{aligned} \tag{52}$$

In equation (52), factors related to \mathbf{Z} are grouped together. We integrated out \mathbf{Z} using

$$\int \exp \left(-\frac{1}{2} (\boldsymbol{\zeta}_j - \boldsymbol{\mu}_j)^\top [\alpha^{-1} \boldsymbol{\Lambda}_j]^{-1} (\boldsymbol{\zeta}_j - \boldsymbol{\mu}_j) \right) d\boldsymbol{\zeta}_j = (2\pi)^{P/2} |\alpha^{-1} \boldsymbol{\Lambda}_j|^{1/2} = (2\pi)^{P/2} \alpha^{-P/2} |\boldsymbol{\Lambda}_j|^{1/2}, \tag{53}$$

which is the normalization constant of P -dimensional Gaussian distribution. Since we have the sum over the number of correspondences (latent vectors), P , the above equation ranges for all of these clusters. The equation (52), becomes

$$p(\Phi|\mathbf{S}, \mathbf{W}, a, b, r) = \frac{b^a}{\Gamma(a)} \int \left(\frac{\alpha}{2\pi} \right)^{\sum_d L_d N_d / 2} \left(\frac{\alpha r}{2\pi} \right)^{PQ/2} \prod_{j=1}^Q \left[(2\pi)^{P/2} \alpha^{-P/2} |\mathbf{\Lambda}_j|^{1/2} \right] \exp \left(-\alpha \left[\frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) - \frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \mathbf{\Lambda}_j^{-1} \boldsymbol{\mu}_j + b \right] \right) \alpha^{a-1} d\alpha \quad (54)$$

$$= \frac{b^a}{\Gamma(a)} \int \left(\frac{\alpha}{2\pi} \right)^{\sum_d L_d N_d / 2} \left(\frac{\alpha r}{2\pi} \right)^{PQ/2} (2\pi)^{PQ/2} \alpha^{-PQ/2} \prod_{j=1}^Q |\mathbf{\Lambda}_j|^{1/2} \exp \left(-\alpha \left[\frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) - \frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \mathbf{\Lambda}_j^{-1} \boldsymbol{\mu}_j + b \right] \right) \alpha^{a-1} d\alpha \quad (55)$$

The α parameter is integrated out by using the following normalization constant of a Gamma distribution

$$\int \alpha^{a'-1} \exp(-b\alpha) d\alpha = \frac{\Gamma(a')}{b^{a'}}. \quad (56)$$

Finally the likelihood of the nonlinear model is given by

$$p(\Phi|\mathbf{S}, \mathbf{W}, a, b, r) = (2\pi)^{-\frac{\sum_d L_d N_d}{2}} r^{\frac{PQ}{2}} \frac{b^a}{b^{a'}} \frac{\Gamma(a')}{\Gamma(a)} \prod_{j=1}^Q |\mathbf{\Lambda}_j|^{1/2}, \quad (57)$$

Here,

$$a' = a + \frac{\sum_d L_d N_d}{2}, \quad b' = b + \frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) - \frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \mathbf{\Lambda}_j^{-1} \boldsymbol{\mu}_j, \quad (58)$$

and

$$\boldsymbol{\mu}_j = \mathbf{\Lambda}_j \sum_{d=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d, \quad \mathbf{\Lambda}_j^{-1} = \sum_{d=1}^D N_{dj} \hat{\mathbf{K}}_d + r \mathbf{I}. \quad (59)$$

1.4 Posterior

The posterior for the precision parameter α is given by

$$p(\alpha|\Phi, \mathbf{S}, \mathbf{W}, a, b) = \mathcal{G}(a', b'), \quad (60)$$

and the posterior for the latent vector $\boldsymbol{\zeta}_j$ is given by

$$p(\boldsymbol{\zeta}_j|\alpha, \Phi, \mathbf{S}, \mathbf{W}, r) = \mathcal{N}(\boldsymbol{\mu}_j, \alpha^{-1} \mathbf{\Lambda}_j) \quad (61)$$

The derivation for these posteriors is given by

$$\begin{aligned}
p(\alpha|\Phi, \mathbf{S}, \mathbf{W}, a, b) \prod_{j=1}^Q p(\zeta_j|\alpha, \Phi, \mathbf{S}, \mathbf{W}, r) &\propto p(\Phi|\alpha, \mathbf{Z}, \mathbf{S}, \mathbf{W}, a, b, r) p(\alpha|a, b) \prod_{j=1}^Q p(\zeta_j|\alpha, r) \\
&= \prod_{d=1}^D \prod_{n=1}^{N_d} \mathcal{N}(\phi(\mathbf{x}_{dn})|\mathbf{W}_d \zeta_{s_{dn}}, \alpha^{-1} \mathbf{I}) \mathcal{G}(\alpha|a, b) \times \prod_{j=1}^Q \mathcal{N}(\zeta_j|\mathbf{0}, (\alpha r)^{-1} \mathbf{I}) \\
&= \prod_{d=1}^D \prod_{n=1}^{N_d} \left(\frac{\alpha}{2\pi}\right)^{L_d/2} \exp\left(-\frac{\alpha}{2} \|\phi(\mathbf{x}_{dn}) - \mathbf{W}_d \zeta_{s_{dn}}\|^2\right) \prod_{j=1}^Q \left(\frac{\alpha r}{2\pi}\right)^{P/2} \\
&\times \exp\left(-\frac{\alpha r}{2} \|\zeta_j\|^2\right) \frac{b^a \alpha^{a-1}}{\Gamma(a)} \exp(-b\alpha) \\
&\propto \alpha^{a'-1} \exp(-b'\alpha) \prod_{j=1}^Q |\Lambda_j|^{-1/2} \exp\left(-\frac{\alpha}{2} (\zeta_j - \mu_j)^\top \Lambda_j^{-1} (\zeta_j - \mu_j)\right) \\
&\propto \mathcal{G}(a', b') \prod_{j=1}^Q \mathcal{N}(\mu_j, \alpha^{-1} \Lambda_j).
\end{aligned} \tag{62}$$

2 Inference for the non-linear model

2.1 M-step

In the M-step, the projection matrices \mathbf{W} are estimated by maximizing the logarithm of the joint likelihood (3). The gradient of the joint likelihood is computed by

$$\frac{\partial \log p(\Phi, \mathbf{S}|\mathbf{W}, a, b, r, \gamma)}{\partial \mathbf{W}_d} = \frac{\partial \log p(\mathbf{S}|\gamma)}{\partial \mathbf{W}_d} + \frac{\partial \log p(\Phi|\mathbf{S}, \mathbf{W}, a, b, r)}{\partial \mathbf{W}_d}, \tag{63}$$

Since the derivative of the first term in the above expression is zero, the expression becomes

$$\begin{aligned}
\frac{\partial \log p(\Phi, \mathbf{S}|\mathbf{W}, a, b, r, \gamma)}{\partial \mathbf{W}_d} &= \frac{\partial \log p(\Phi|\mathbf{S}, \mathbf{W}, a, b, r)}{\partial \mathbf{W}_d} = \frac{\partial \log \left[(2\pi)^{-\frac{\sum_d L_d N_d}{2}} r^{\frac{PQ}{2}} \frac{b^a}{b'^{a'}} \frac{\Gamma(a')}{\Gamma(a)} \prod_{j=1}^Q |\Lambda_j|^{1/2} \right]}{\partial \mathbf{W}_d}, \\
&= \frac{\partial \log \left[\frac{cte}{b'^{a'}} \prod_{j=1}^Q |\Lambda_j|^{1/2} \right]}{\partial \mathbf{W}_d}, \\
&= -\frac{a'}{b'} \frac{\partial b'}{\partial \mathbf{W}_d} + \frac{1}{2} \sum_{j=1}^Q \frac{\partial \log |\Lambda_j|}{\partial \mathbf{W}_d} + 0, \\
&= -\frac{a'}{b'} \frac{\partial b'}{\partial \mathbf{W}_d} + \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\Lambda_j^{-1} \frac{\partial \Lambda_j}{\partial \mathbf{W}_d} \right)
\end{aligned} \tag{64}$$

where

$$\frac{\partial b'}{\partial \mathbf{W}_d} = \frac{\partial}{\partial \mathbf{W}_d} \left[\frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) - \frac{1}{2} \sum_{j=1}^Q \mu_j^\top \Lambda_j^{-1} \mu_j + b \right], \tag{65}$$

Here the second factor of the argument is the only which depends on \mathbf{W}_d .

$$\begin{aligned}
\frac{\partial b'}{\partial \mathbf{W}_d} &= \frac{\partial}{\partial \mathbf{W}_d} \left[-\frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right] \\
&= -\frac{1}{2} \sum_{j=1}^Q \frac{\partial}{\partial \mathbf{W}_d} \left[\boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right] \\
&= -\frac{1}{2} \sum_{j=1}^Q \frac{\partial}{\partial \mathbf{W}_d} \left[\left(\boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d \right)^\top \boldsymbol{\Lambda}_j^{-1} \left(\boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d \right) \right] \\
&= -\frac{1}{2} \sum_{j=1}^Q \frac{\partial}{\partial \mathbf{W}_d} \left[\sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d^\top \boldsymbol{\Lambda}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d \right] \\
&= -\frac{1}{2} \sum_{j=1}^Q \frac{\partial}{\partial \mathbf{W}_d} \left[\sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d^\top \boldsymbol{\Lambda}_j^\top \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d \right] \tag{66}
\end{aligned}$$

Since the derivative with respect to \mathbf{W}_d takes the d th-projection matrix, the sum over $e = 1, \dots, d, \dots, D$ only affect the d th-projection matrix. Besides, by using product derivatives properties, we must to compute the derivative of

$$\frac{\partial b'}{\partial \mathbf{W}_d} = -\frac{1}{2} \left[\sum_{n:s_{dn}=j} \left(\partial \hat{\mathbf{k}}_d \right)^\top \boldsymbol{\Lambda}_j^\top \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d + \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d^\top \left[\left(\partial \boldsymbol{\Lambda}_j^\top \right) \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d + \boldsymbol{\Lambda}_j^\top \sum_{n:s_{dn}=j} \partial \hat{\mathbf{k}}_d \right] \right] \tag{67}$$

Since $\hat{\mathbf{k}}_d = \mathbf{W}_d^\top \boldsymbol{\phi}(\mathbf{x}_{dn})$ the above derivative $\partial \hat{\mathbf{k}}_d$ becomes,

$$\frac{\partial \hat{\mathbf{k}}_d}{\partial \mathbf{W}_d} = \frac{\partial}{\partial \mathbf{W}_d} \left[\mathbf{W}_d^\top \boldsymbol{\phi}(\mathbf{x}_{dn}) \right] = \boldsymbol{\phi}(\mathbf{x}_{dn}). \tag{68}$$

For the $\boldsymbol{\Lambda}_j$, the derivative is computed as

$$\begin{aligned}
\frac{\partial \boldsymbol{\Lambda}_j}{\partial \mathbf{W}_d} &= \frac{\partial}{\partial \mathbf{W}_d} \left(\left[\sum_{e=1}^D N_{dj} \hat{\mathbf{K}}_d + r \mathbf{I} \right]^{-1} \right), \\
&= -\boldsymbol{\Lambda}_j \frac{\partial}{\partial \mathbf{W}_d} \left(\sum_{e=1}^D N_{dj} \hat{\mathbf{K}}_d + r \mathbf{I} \right) \boldsymbol{\Lambda}_j, \\
&= -N_{dj} \boldsymbol{\Lambda}_j \frac{\partial \hat{\mathbf{K}}_d}{\partial \mathbf{W}_d} \boldsymbol{\Lambda}_j, \tag{69}
\end{aligned}$$

where $\hat{\mathbf{K}}_d = \mathbf{W}_d^\top \mathbf{W}_d$. The derivative $\frac{\partial \hat{\mathbf{K}}_d}{\partial \mathbf{W}_d}$ is given by

$$\frac{\partial \hat{\mathbf{K}}_d}{\partial \mathbf{W}_d} = \frac{\partial}{\partial \mathbf{W}_d} \left[\mathbf{W}_d^\top \mathbf{W}_d \right] = 2\mathbf{W}_d \tag{70}$$

By replacing (67) and (69) in (64), the gradient for \mathbf{W}_d becomes

$$\begin{aligned}
\frac{\partial \log p(\Phi, \mathbf{S} | \mathcal{W}, a, b, r, \gamma)}{\partial \mathcal{W}_d} &= -\frac{a'}{b'} \frac{\partial b'}{\partial \mathcal{W}_d} + \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\Lambda_j^{-1} \frac{\partial \Lambda_j}{\partial \mathcal{W}_d} \right), \\
&= -\frac{a'}{b'} \frac{\partial b'}{\partial \mathcal{W}_d} - \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\Lambda_j^{-1} N_{dj} \Lambda_j \frac{\partial \hat{\mathbf{K}}_d}{\partial \mathcal{W}_d} \Lambda_j \right), \\
&= -\frac{a'}{b'} \frac{\partial b'}{\partial \mathcal{W}_d} - \frac{1}{2} \sum_{j=1}^Q N_{dj} \text{Tr} \left(\frac{\partial \hat{\mathbf{K}}_d}{\partial \mathcal{W}_d} \Lambda_j \right), \\
&= \frac{1}{2} \frac{a'}{b'} \left[\sum_{n:s_{dn}=j} \left(\partial \hat{\mathbf{k}}_d \right)^\top \Lambda_j^\top \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d \right. \\
&\quad \left. + \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d^\top \left[\left(\partial \Lambda_j^\top \right) \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d + \Lambda_j^\top \sum_{n:s_{dn}=j} \partial \hat{\mathbf{k}}_d \right] \right] \\
&\quad - \frac{1}{2} \sum_{j=1}^Q N_{dj} \text{Tr} \left(\frac{\partial \hat{\mathbf{K}}_d}{\partial \mathcal{W}_d} \Lambda_j \right), \\
&= \frac{1}{2} \frac{a'}{b'} \left[\sum_{n:s_{dn}=j} \left(\partial \hat{\mathbf{k}}_d \right)^\top \boldsymbol{\mu}_j + \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d^\top \right. \\
&\quad \left. \times \left[\left(\partial \Lambda_j^\top \right) \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d + \Lambda_j^\top \sum_{n:s_{dn}=j} \partial \hat{\mathbf{k}}_d \right] \right] \\
&\quad - \frac{1}{2} \sum_{j=1}^Q N_{dj} \text{Tr} \left(\frac{\partial \hat{\mathbf{K}}_d}{\partial \mathcal{W}_d} \Lambda_j \right), \\
&= \frac{1}{2} \frac{a'}{b'} \left[\sum_{n:s_{dn}=j} \phi(\mathbf{x}_{dn})^\top \boldsymbol{\mu}_j + \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d^\top \right. \\
&\quad \left. \times \left[\left(-N_{dj} \Lambda_j \frac{\partial \hat{\mathbf{K}}_d}{\partial \mathcal{W}_d} \Lambda_j \right)^\top \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d + \Lambda_j^\top \sum_{n:s_{dn}=j} \partial \hat{\mathbf{k}}_d \right] \right] \\
&\quad - \frac{1}{2} \sum_{j=1}^Q N_{dj} \text{Tr} \left(\frac{\partial \hat{\mathbf{K}}_d}{\partial \mathcal{W}_d} \Lambda_j \right), \\
&= \frac{1}{2} \frac{a'}{b'} \left[\sum_{n:s_{dn}=j} \phi(\mathbf{x}_{dn})^\top \boldsymbol{\mu}_j - 2N_{dj} \boldsymbol{\mu}_j^\top \mathcal{W}_d \boldsymbol{\mu}_j + \boldsymbol{\mu}_j^\top \sum_{n:s_{dn}=j} \phi(\mathbf{x}_{dn}) \right] \\
&\quad - \frac{1}{2} \sum_{j=1}^Q 2N_{dj} \text{Tr}(\mathcal{W}_d \Lambda_j),
\end{aligned} \tag{71}$$

$$\begin{aligned}
&= \frac{1}{2} \frac{a'}{b'} \left[\sum_{n:s_{dn}=j} \phi(\mathbf{x}_{dn})^\top \boldsymbol{\mu}_j - 2N_{dj} \boldsymbol{\mu}_j^\top \mathcal{W}_d \boldsymbol{\mu}_j + \boldsymbol{\mu}_j^\top \sum_{n:s_{dn}=j} \phi(\mathbf{x}_{dn}) \right] \\
&\quad - \frac{1}{2} \sum_{j=1}^Q 2N_{dj} \text{Tr}(\mathcal{W}_d \Lambda_j),
\end{aligned} \tag{72}$$

2.2 Kernel derivatives

Since $\hat{\mathbf{k}}_d$ depends on \mathcal{W}_d ,

$$\frac{\partial \hat{\mathbf{k}}_d}{\partial \mathcal{W}_d} = \frac{\partial}{\partial \mathcal{W}_d} \begin{bmatrix} k(\mathbf{w}_1^d, \mathbf{x}_{dn}) \\ k(\mathbf{w}_2^d, \mathbf{x}_{dn}) \\ \vdots \\ k(\mathbf{w}_P^d, \mathbf{x}_{dn}) \end{bmatrix} = \frac{\partial \hat{\mathbf{k}}_d}{\partial \boldsymbol{\Theta}} \tag{73}$$

where $\boldsymbol{\Theta} = \{l, \sigma^2, \{\mathbf{w}_i\}_{i=1}^P\}$ are the kernel hyperparameters of a squared exponential kernel of the form $k(x, x') = \sigma^2 \exp\left(-\frac{(x-x')^2}{2l^2}\right)$.

$$\frac{\partial \hat{\mathbf{K}}_d}{\partial \Theta} = \begin{bmatrix} \frac{\partial k(\mathbf{w}_1^d, \mathbf{x}_{dn})}{\partial \sigma^2} + \frac{\partial k(\mathbf{w}_2^d, \mathbf{x}_{dn})}{\partial \sigma^2} + \dots + \frac{\partial k(\mathbf{w}_P^d, \mathbf{x}_{dn})}{\partial \sigma^2} \\ \frac{\partial k(\mathbf{w}_1^d, \mathbf{x}_{dn})}{\partial l^2} + \frac{\partial k(\mathbf{w}_2^d, \mathbf{x}_{dn})}{\partial l^2} + \dots + \frac{\partial k(\mathbf{w}_P^d, \mathbf{x}_{dn})}{\partial l^2} \\ \frac{\partial k(\mathbf{w}_1^d, \mathbf{x}_{dn})}{\partial \mathbf{w}_1^d} + \frac{\partial k(\mathbf{w}_2^d, \mathbf{x}_{dn})}{\partial \mathbf{w}_1^d} + \dots + \frac{\partial k(\mathbf{w}_P^d, \mathbf{x}_{dn})}{\partial \mathbf{w}_1^d} \\ \frac{\partial k(\mathbf{w}_1^d, \mathbf{x}_{dn})}{\partial \mathbf{w}_P^d} + \frac{\partial k(\mathbf{w}_2^d, \mathbf{x}_{dn})}{\partial \mathbf{w}_P^d} + \dots + \frac{\partial k(\mathbf{w}_P^d, \mathbf{x}_{dn})}{\partial \mathbf{w}_P^d} \end{bmatrix}^\top \quad (74)$$

Let us define the following derivatives with respect to the kernel hyperparameters in the equation (74), $k(\mathbf{w}_i^d, \mathbf{x}_{dn}) = \sigma^2 \exp\left(-\frac{\beta}{2} (\mathbf{w}_i^d - \mathbf{x}_{dn})^\top (\mathbf{w}_i^d - \mathbf{x}_{dn})\right)$, where β is the inverwidth $\beta = 1/l^2$

$$\frac{\partial k(\mathbf{w}_i^d, \mathbf{x}_{dn})}{\partial \sigma^2} = \frac{1}{\sigma^2} k(\mathbf{w}_i^d, \mathbf{x}_{dn}), \quad (75)$$

$$\frac{\partial k(\mathbf{w}_i^d, \mathbf{x}_{dn})}{\partial \beta} = -(1/2) (\mathbf{w}_i^d - \mathbf{x}_{dn})^\top (\mathbf{w}_i^d - \mathbf{x}_{dn}) k(\mathbf{w}_i^d, \mathbf{x}_{dn}), \quad (76)$$

$$\frac{\partial k(\mathbf{w}_i^d, \mathbf{x}_{dn})}{\partial \mathbf{w}_i^d} = -\beta (\mathbf{w}_i^d - \mathbf{x}_{dn}) k(\mathbf{w}_i^d, \mathbf{x}_{dn}) \quad (77)$$

and

$$\frac{\partial k(\mathbf{w}_i^d, \mathbf{x}_{dn})}{\partial \mathbf{w}_j^d} = 0, \quad \forall i \neq j. \quad (78)$$

Besides for the derivative of $\frac{\partial \hat{\mathbf{K}}_d}{\partial \mathbf{W}_d}$, we have,

$$\frac{\partial \hat{\mathbf{K}}_d}{\partial \mathbf{W}_d} = \begin{bmatrix} \frac{\partial k(\mathbf{w}_1^d, \mathbf{w}_1^d)}{\partial \Theta} & \frac{\partial k(\mathbf{w}_1^d, \mathbf{w}_2^d)}{\partial \Theta} & \dots & \frac{\partial k(\mathbf{w}_1^d, \mathbf{w}_P^d)}{\partial \Theta} \\ \frac{\partial k(\mathbf{w}_2^d, \mathbf{w}_1^d)}{\partial \Theta} & \frac{\partial k(\mathbf{w}_2^d, \mathbf{w}_2^d)}{\partial \Theta} & \dots & \frac{\partial k(\mathbf{w}_2^d, \mathbf{w}_P^d)}{\partial \Theta} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial k(\mathbf{w}_P^d, \mathbf{w}_1^d)}{\partial \Theta} & \frac{\partial k(\mathbf{w}_P^d, \mathbf{w}_2^d)}{\partial \Theta} & \dots & \frac{\partial k(\mathbf{w}_P^d, \mathbf{w}_P^d)}{\partial \Theta} \end{bmatrix} \quad (79)$$

3 Derivatives for the optimization problem

In the M-step, the projection matrices \mathbf{W} are estimated by maximizing the logarithm of the joint likelihood (3). The gradient of the joint likelihood is computed by

$$\frac{\partial \log p(\Phi, \mathbf{S} | \mathbf{W}, a, b, r, \gamma)}{\partial \mathbf{W}_d} = \frac{\partial \log p(\mathbf{S} | \gamma)}{\partial \mathbf{W}_d} + \frac{\partial \log p(\Phi | \mathbf{S}, \mathbf{W}, a, b, r)}{\partial \mathbf{W}_d}, \quad (80)$$

Since the derivative of the first term in the above expression is zero, the expression becomes

$$\begin{aligned} \frac{\partial \log p(\Phi, \mathbf{S} | \mathbf{W}, a, b, r, \gamma)}{\partial \mathbf{W}_d} &= \frac{\partial \log p(\Phi | \mathbf{S}, \mathbf{W}, a, b, r)}{\partial \mathbf{W}_d} = \frac{\partial \log \left[(2\pi)^{-\frac{\sum_d L_d N_d}{2}} r^{\frac{PQ}{2}} \frac{b^a}{b'^a} \frac{\Gamma(a')}{\Gamma(a)} \prod_{j=1}^Q |\Lambda_j|^{1/2} \right]}{\partial \mathbf{W}_d}, \\ &= \frac{\partial \log \left[\frac{cte}{b'^a} \prod_{j=1}^Q |\Lambda_j|^{1/2} \right]}{\partial \mathbf{W}_d}, \\ &= -\frac{a'}{b'} \frac{\partial b'}{\partial \mathbf{W}_d} + \frac{1}{2} \sum_{j=1}^Q \frac{\partial \log |\Lambda_j|}{\partial \mathbf{W}_d} + 0, \\ &= -\frac{a'}{b'} \frac{\partial b'}{\partial \mathbf{W}_d} + \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\Lambda_j^{-1} \frac{\partial \Lambda_j}{\partial \mathbf{W}_d} \right) \end{aligned} \quad (81)$$

We set three different kernels to analyze $k(\mathbf{x}_{dn}, \mathbf{x}_{dn})$, $k(\mathbf{w}_i^d, \mathbf{x}_{dn})$ and $k(\mathbf{w}_i^d, \mathbf{w}_j^d)$. Then, by using a rbf kernel for each one, the hyperparameters are set as:

$$\boldsymbol{\theta} = \left\{ \beta_{xx}, \quad \sigma_{xx}^2, \quad \beta_{xw}, \quad \sigma_{xw}^2, \quad \beta_{ww}, \quad \sigma_{ww}^2, \quad \{\mathbf{w}_i\}_{i=1}^P \right\} \quad (82)$$

First, for $k(\mathbf{x}_{dn}, \mathbf{x}_{dn})$ we have the following derivatives,

$$\frac{\partial \log p(\boldsymbol{\Phi}, \mathbf{S} | \boldsymbol{\mathcal{W}}, a, b, r, \gamma)}{\partial \beta_{xx}} = -\frac{a'}{b'} \frac{\partial b'}{\partial \beta_{xx}} \quad (83)$$

$$= -\frac{a'}{b'} \frac{\partial}{\partial \beta_{xx}} \left\{ b + \frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) - \frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right\} \quad (84)$$

$$= -\frac{a'}{b'} \left\{ \frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \frac{\partial}{\partial \beta_{xx}} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) \right\} \quad (85)$$

$$= \frac{1}{4} \frac{a'}{b'} \sum_{d=1}^D \sum_{n=1}^{N_d} (\mathbf{x}_{dn} - \mathbf{x}_{dn})^\top (\mathbf{x}_{dn} - \mathbf{x}_{dn}) k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) \quad (86)$$

$$= \mathbf{0}, \quad (87)$$

for the kernel inversewidth β_{xx} ,

$$\frac{\partial \log p(\boldsymbol{\Phi}, \mathbf{S} | \boldsymbol{\mathcal{W}}, a, b, r, \gamma)}{\partial \sigma_{xx}^2} = -\frac{a'}{b'} \frac{\partial b'}{\partial \sigma_{xx}^2} \quad (88)$$

$$= -\frac{a'}{b'} \frac{\partial}{\partial \sigma_{xx}^2} \left\{ b + \frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) - \frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right\} \quad (89)$$

$$= -\frac{a'}{b'} \left\{ \frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} \frac{\partial}{\partial \sigma_{xx}^2} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) \right\} \quad (90)$$

$$= -\frac{a'}{2b'} \sum_{d=1}^D \sum_{n=1}^{N_d} \frac{1}{\sigma_{xx}^2} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) \quad (91)$$

for the kernel variance σ_{xx}^2 .

Besides the kernel derivatives for the $k(\mathbf{w}_i^d, \mathbf{x}_{dn})$ are

$$\frac{\partial \log p(\Phi, \mathbf{S} | \mathcal{W}, a, b, r, \gamma)}{\partial \beta_{wx}} = -\frac{a'}{b'} \frac{\partial b'}{\partial \beta_{wx}}, \quad (92)$$

$$= -\frac{a'}{b'} \frac{\partial}{\partial \beta_{wx}} \left\{ b + \frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) - \frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right\}, \quad (93)$$

$$= -\frac{a'}{b'} \frac{\partial}{\partial \beta_{wx}} \left\{ -\frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right\}, \quad (94)$$

$$= -\frac{a'}{b'} \frac{\partial}{\partial \beta_{wx}} \left\{ \left(\boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d \right)^\top \boldsymbol{\Lambda}_j^{-1} \left(\boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d \right) \right\}, \quad (95)$$

$$= \frac{a'}{2b'} \left\{ \left(\boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} \partial \hat{\mathbf{k}}_d \right)^\top \boldsymbol{\Lambda}_j^{-1} \left(\boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d \right) \right. \quad (96)$$

$$\left. + \left(\boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} \hat{\mathbf{k}}_d \right)^\top \boldsymbol{\Lambda}_j^{-1} \left(\boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} \partial \hat{\mathbf{k}}_d \right) \right\}, \quad (97)$$

$$= \frac{a'}{2b'} \left\{ \left(\boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} -\frac{1}{2} (\mathbf{w}_i^d - \mathbf{x}_{dn})^\top (\mathbf{w}_i^d - \mathbf{x}_{dn}) \hat{\mathbf{k}}_d \right)^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right. \quad (98)$$

$$\left. + \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \left(\boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} -\frac{1}{2} (\mathbf{w}_i^d - \mathbf{x}_{dn})^\top (\mathbf{w}_i^d - \mathbf{x}_{dn}) \hat{\mathbf{k}}_d \right) \right\}, \quad (99)$$

for the β_{wx} parameter,

$$\frac{\partial \log p(\Phi, \mathbf{S} | \mathcal{W}, a, b, r, \gamma)}{\partial \sigma_{wx}^2} = -\frac{a'}{b'} \frac{\partial b'}{\partial \sigma_{wx}^2}, \quad (100)$$

$$= -\frac{a'}{b'} \frac{\partial}{\partial \sigma_{wx}^2} \left\{ b + \frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) - \frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right\}, \quad (101)$$

$$= -\frac{a'}{b'} \frac{\partial}{\partial \sigma_{wx}^2} \left\{ -\frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right\}, \quad (102)$$

$$= \frac{a'}{2b'} \left\{ \left(\boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} \frac{1}{\sigma_{wx}^2} \hat{\mathbf{k}}_d \right)^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right. \quad (103)$$

$$\left. + \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \left(\boldsymbol{\Lambda}_j \sum_{e=1}^D \sum_{n:s_{dn}=j} \frac{1}{\sigma_{wx}^2} \hat{\mathbf{k}}_d \right) \right\}, \quad (104)$$

for the σ_{wx} parameter,

$$\frac{\partial \log p(\Phi, \mathbf{S} | \mathbf{W}, a, b, r, \gamma)}{\partial \mathbf{w}_i^d} = -\frac{a'}{b'} \frac{\partial b'}{\partial \mathbf{w}_i^d}, \quad (105)$$

$$= -\frac{a'}{b'} \frac{\partial}{\partial \mathbf{w}_i^d} \left\{ b + \frac{1}{2} \sum_{d=1}^D \sum_{n=1}^{N_d} k(\mathbf{x}_{dn}, \mathbf{x}_{dn}) - \frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right\}, \quad (106)$$

$$= -\frac{a'}{b'} \frac{\partial}{\partial \mathbf{w}_i^d} \left\{ -\frac{1}{2} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right\}, \quad (107)$$

$$= \frac{a'}{2b'} \left\{ \left(\boldsymbol{\Lambda}_j \sum_{n:s_{dn}=j} -\beta_{wx} (\mathbf{w}_i^d - \mathbf{x}_{dn}) \hat{\mathbf{k}}_d \right)^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right. \quad (108)$$

$$\left. + \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \left(\boldsymbol{\Lambda}_j \sum_{n:s_{dn}=j} -\beta_{wx} (\mathbf{w}_i^d - \mathbf{x}_{dn}) \hat{\mathbf{k}}_d \right) \right\}, \quad (109)$$

for the \mathbf{w}_i^d parameter.

The last kernel is defined as $k(\mathbf{w}_i^d, \mathbf{w}_{i'}^d)$, with hyperparameters σ_{ww}^2 , β_{ww} and \mathbf{w}_i^d . These derivatives are computed as follows

$$\frac{\partial \log p(\Phi, \mathbf{S} | \mathbf{W}, a, b, r, \gamma)}{\partial \beta_{ww}} = \left[-\frac{a'}{b'} \frac{\partial b'}{\partial \beta_{ww}} + \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\boldsymbol{\Lambda}_j^{-1} \frac{\partial \boldsymbol{\Lambda}_j}{\partial \beta_{ww}} \right) \right], \quad (110)$$

$$= \frac{\partial}{\partial \beta_{ww}} \left\{ \frac{a'}{2b'} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \boldsymbol{\Lambda}_j^{-1} \boldsymbol{\mu}_j \right\} + \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\boldsymbol{\Lambda}_j^{-1} \frac{\partial \boldsymbol{\Lambda}_j}{\partial \beta_{ww}} \right), \quad (111)$$

$$= \frac{\partial}{\partial \beta_{ww}} \left\{ \frac{a'}{2b'} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \left(\sum_{d=1}^D N_{dj} \hat{\mathbf{K}}_d + r \mathbf{I} \right) \boldsymbol{\mu}_j \right\} \quad (112)$$

$$+ \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\boldsymbol{\Lambda}_j^{-1} \frac{\partial}{\partial \beta_{ww}} \left(\sum_{d=1}^D N_{dj} \hat{\mathbf{K}}_d + r \mathbf{I} \right)^{-1} \right), \quad (113)$$

$$= \frac{\partial}{\partial \beta_{ww}} \left\{ \frac{a'}{2b'} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \left(\sum_{d=1}^D N_{dj} \hat{\mathbf{K}}_d + r \mathbf{I} \right) \boldsymbol{\mu}_j \right\} \quad (114)$$

$$+ \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\boldsymbol{\Lambda}_j^{-1} \frac{\partial}{\partial \beta_{ww}} \left(\sum_{d=1}^D N_{dj} \hat{\mathbf{K}}_d + r \mathbf{I} \right)^{-1} \right), \quad (115)$$

$$= \left\{ \frac{a'}{2b'} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \frac{\partial}{\partial \beta_{ww}} \left(\sum_{d=1}^D N_{dj} \hat{\mathbf{K}}_d \right) \boldsymbol{\mu}_j \right\} \quad (116)$$

$$- \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\boldsymbol{\Lambda}_j^{-1} \boldsymbol{\Lambda}_j \frac{\partial}{\partial \beta_{ww}} \left(\sum_{d=1}^D N_{dj} \hat{\mathbf{K}}_d + r \mathbf{I} \right) \boldsymbol{\Lambda}_j \right), \quad (117)$$

$$= \left\{ \frac{a'}{2b'} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \left(\sum_{d=1}^D N_{dj} \left(-(1/2) (\mathbf{P}_2) \hat{\mathbf{K}}_d \right) \right) \boldsymbol{\mu}_j \right\} \quad (118)$$

$$- \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\left(\sum_{d=1}^D N_{dj} \left(-(1/2) (\mathbf{P}_2) \right) \hat{\mathbf{K}}_d \right) \boldsymbol{\Lambda}_j \right), \quad (119)$$

where \mathbf{P}_2 is the pairwise distance of \mathbf{W}_d

$$\frac{\partial \log p(\Phi, \mathbf{S} | \mathbf{W}, a, b, r, \gamma)}{\partial \sigma_{ww}^2} = \left[-\frac{a'}{b'} \frac{\partial b'}{\partial \sigma_{ww}^2} + \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\Lambda_j^{-1} \frac{\partial \Lambda_j}{\partial \sigma_{ww}^2} \right) \right], \quad (120)$$

$$= \left\{ \frac{a'}{2b'} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \frac{\partial}{\partial \sigma_{ww}^2} \left(\sum_{d=1}^D N_{dj} \hat{\mathbf{K}}_d \right) \boldsymbol{\mu}_j \right\} \quad (121)$$

$$- \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\Lambda_j^{-1} \Lambda_j \frac{\partial}{\partial \sigma_{ww}^2} \left(\sum_{d=1}^D N_{dj} \hat{\mathbf{K}}_d + r \mathbf{I} \right) \Lambda_j \right), \quad (122)$$

$$= \left\{ \frac{a'}{2b'} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \left(\sum_{d=1}^D N_{dj} \left(\frac{1}{\sigma_{ww}^2} \hat{\mathbf{K}}_d \right) \right) \boldsymbol{\mu}_j \right\} \quad (123)$$

$$- \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\left(\sum_{d=1}^D N_{dj} \frac{1}{\sigma_{ww}^2} \hat{\mathbf{K}}_d \right) \Lambda_j \right). \quad (124)$$

Finally for the \mathbf{w}_i^d parameter we have

$$\frac{\partial \log p(\Phi, \mathbf{S} | \mathbf{W}, a, b, r, \gamma)}{\partial \mathbf{w}_i^d} = \left[-\frac{a'}{b'} \frac{\partial b'}{\partial \mathbf{w}_i^d} + \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\Lambda_j^{-1} \frac{\partial \Lambda_j}{\partial \mathbf{w}_i^d} \right) \right], \quad (125)$$

$$= \left\{ \frac{a'}{2b'} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \frac{\partial}{\partial \mathbf{w}_i^d} \left(\sum_{d=1}^D N_{dj} \hat{\mathbf{K}}_d \right) \boldsymbol{\mu}_j \right\} \quad (126)$$

$$- \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\Lambda_j^{-1} \Lambda_j \frac{\partial}{\partial \mathbf{w}_i^d} \left(\sum_{d=1}^D N_{dj} \hat{\mathbf{K}}_d + r \mathbf{I} \right) \Lambda_j \right), \quad (127)$$

$$= \left\{ \frac{a'}{2b'} \sum_{j=1}^Q \boldsymbol{\mu}_j^\top \left(\sum_{d=1}^D N_{dj} \left(-\beta_{ww} \mathbf{P}_1 \hat{\mathbf{K}}_d \right) \right) \boldsymbol{\mu}_j \right\} \quad (128)$$

$$- \frac{1}{2} \sum_{j=1}^Q \text{Tr} \left(\left(\sum_{d=1}^D N_{dj} - \beta_{ww} \mathbf{P}_1 \hat{\mathbf{K}}_d \right) \Lambda_j \right), \quad (129)$$

where \mathbf{P}_1 define mathematical notation

References

- [1] Tomoharu Iwata, Tsutomu Hirao, and Naonori Ueda. Unsupervised cluster matching via probabilistic latent variable models, 2013.
- [2] Tomoharu Iwata, Tsutomu Hirao, and Naonori Ueda. Probabilistic latent variable models for unsupervised many-to-many object matching. *Information Processing and Management*, 52(4):682 – 697, 2016.