



UNIVERSITY OF AMSTERDAM

University of Amsterdam

# Homework Assignment 3

Machine Learning I

2024

Pedro M. P. Curvo

15713725

# Contents

1	Principal Component Analysis . . . . .	1
2	Probabilistic PCA - A general latent space distribution . . . . .	11
3	Mixtures of Experts . . . . .	13

# 1 Principal Component Analysis

a)

i)

Starting from the formalization of maximizing the scatter we have that:

$$\max \sum_{i=1}^n \|P\mathbf{x}_i - P\bar{\mathbf{x}}\|^2$$

Expanding the norm we have that:

$$\begin{aligned} \max \sum_{i=1}^n \|P\mathbf{x}_i - P\bar{\mathbf{x}}\|^2 &= \max \sum_{i=1}^n (P\mathbf{x}_i - P\bar{\mathbf{x}})^T (P\mathbf{x}_i - P\bar{\mathbf{x}}) \\ &= \max \sum_{i=1}^n (P(\mathbf{x}_i - \bar{\mathbf{x}}))^T P(\mathbf{x}_i - \bar{\mathbf{x}}) \\ &= \max \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^T P^T P(\mathbf{x}_i - \bar{\mathbf{x}}) \end{aligned}$$

Now, we need to consider two properties of the matrix  $P$ :

- $P$  is idempotent, i.e.  $P^n = P \quad \forall n \in \mathbb{N}$ , because projecting a vector twice is the same as projecting it once. If it is already in the subspace of the projection than another projection will not change it.
- $P$  is symmetric, i.e.  $P = P^T$

Now, we can rewrite the expression as:

$$\max \sum_{i=1}^n \|P\mathbf{x}_i - P\bar{\mathbf{x}}\|^2 = \max \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^T P(\mathbf{x}_i - \bar{\mathbf{x}})$$

Because,  $P^T P = P P = P$ . This is the same as the expression we wanted to show.

ii)

To prove that it is the same as  $\max \text{Tr}(\mathbf{S}_1 P)$ , we first have to consider the Scatter matrix  $\mathbf{S}_1$ :

$$\mathbf{S}_1 = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

This matrix is often used to estimate the covariance matrix and actually measures the scatter of the data by taking the outer product of the data points.

Now, we can rewrite the expression from the previous question as:

$$\begin{aligned}
\max \sum_{i=1}^n \|P\mathbf{x}_i - P\bar{\mathbf{x}}\|^2 &= \max \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^T P (\mathbf{x}_i - \bar{\mathbf{x}}) \\
&= \max \sum_{i=1}^n \text{Tr} \left( (\mathbf{x}_i - \bar{\mathbf{x}})^T P (\mathbf{x}_i - \bar{\mathbf{x}}) \right) \quad \text{because the trace of a scalar is the scalar itself}
\end{aligned}$$

Because the trace is invariant under cyclic permutations, i.e.  $\text{Tr}(ABC) = \text{Tr}(BCA) = \text{Tr}(CAB)$ , we can rewrite the expression as:

$$\begin{aligned}
&= \max \sum_{i=1}^n \text{Tr} \left( P (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \right) \\
&= \max \text{Tr} \left( P \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \right) \\
&= \max \text{Tr} (P\mathbf{S}_1) \quad \text{because } \mathbf{S}_1 = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T
\end{aligned}$$

This is the same as the expression we wanted to show.

**b)**

**i)**

Centering the data, which is done by subtracting the mean from each data point, is important for PCA because it removes the bias from the data. i.e., the bias introduced by the mean. If the data is not centered, the first principal component may align more with the mean rather than with the covariance, leading to errors.

Besides that, when the data is not centered, the reconstruction error will be higher. This happens because the distance between the original data points and their projections onto the principal components will increase. The projection space is centered at the origin, while the original data may not be. This results in a constant bias in the distances, which cannot be corrected by the projection. Therefore, the reconstruction error needs to account for this bias, complicating the minimization process, specially because the bias is not constant in all directions, e.g, the bias in the direction of the first principal component might be twice as large as the bias in the direction of the second principal component. However, when the data is centered, the bias is removed and this bias distance is 0 in all directions, facilitating the minimization process and not introducing possible errors.

**ii)**

The reconstruction error is given, after centering, by:

$$\sum_{i=1}^n \|(\mathbf{x}_i - \bar{\mathbf{x}}) - \mathbf{P}(\mathbf{x}_i - \bar{\mathbf{x}})\|^2$$

If we isolate  $(\mathbf{x}_i - \bar{\mathbf{x}})$  and expand the norm, we have:

$$\begin{aligned} \sum_{i=1}^n \|(\mathbf{x}_i - \bar{\mathbf{x}}) - \mathbf{P}(\mathbf{x}_i - \bar{\mathbf{x}})\|^2 &= \sum_{i=1}^n \|(\mathbf{I} - \mathbf{P})(\mathbf{x}_i - \bar{\mathbf{x}})\|^2 \\ &= \sum_{i=1}^n ((\mathbf{I} - \mathbf{P})(\mathbf{x}_i - \bar{\mathbf{x}}))^T (\mathbf{I} - \mathbf{P})(\mathbf{x}_i - \bar{\mathbf{x}}) \\ &= \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^T (\mathbf{I} - \mathbf{P})^T (\mathbf{I} - \mathbf{P})(\mathbf{x}_i - \bar{\mathbf{x}}) \end{aligned}$$

Now, since  $\mathbf{P}$  is idempotent and symmetric, we have that  $(\mathbf{I} - \mathbf{P})^T (\mathbf{I} - \mathbf{P}) = (\mathbf{I} - \mathbf{P})(\mathbf{I} - \mathbf{P}) = \mathbf{I} - \mathbf{P} - \mathbf{P} + \mathbf{P}^2 = \mathbf{I} - \mathbf{P}$ .

Therefore, the expression becomes:

$$\min \sum_{i=1}^n \|(\mathbf{x}_i - \bar{\mathbf{x}}) - \mathbf{P}(\mathbf{x}_i - \bar{\mathbf{x}})\|^2 = \min \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^T (\mathbf{I} - \mathbf{P})(\mathbf{x}_i - \bar{\mathbf{x}})$$

As we wanted to show.

iii)

Similar to the previous question (a)ii),  $S_2$  is defined as:

$$\mathbf{S}_2 = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T$$

And is the scatter matrix.

Following the same steps as before, we have that:

$$\begin{aligned} \min \sum_{i=1}^n \|(\mathbf{x}_i - \bar{\mathbf{x}}) - \mathbf{P}(\mathbf{x}_i - \bar{\mathbf{x}})\|^2 &= \min \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^T (\mathbf{I} - \mathbf{P}) (\mathbf{x}_i - \bar{\mathbf{x}}) \\ &= \min \sum_{i=1}^n \text{Tr} \left( (\mathbf{x}_i - \bar{\mathbf{x}})^T (\mathbf{I} - \mathbf{P}) (\mathbf{x}_i - \bar{\mathbf{x}}) \right) \\ &= \min \text{Tr} \left( (\mathbf{I} - \mathbf{P}) \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \right) \\ &= \min \text{Tr} (\mathbf{S}_2 (\mathbf{I} - \mathbf{P})) \quad \text{because } \mathbf{S}_2 = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \end{aligned}$$

As we wanted to show.

c)

i)

Intuitively, a projection of a vector cannot have a higher length than the original vector, hence  $\|P\mathbf{y}\|^2 \leq \|\mathbf{y}\|^2$ .

We can prove it by decomposing the vector  $\mathbf{y}$  into two components, one that is in the subspace of the projection and another that is orthogonal to it.

$$\mathbf{y} = P\mathbf{y} + (I - P)\mathbf{y}$$

Meaning that:

$$\begin{aligned}\|\mathbf{y}\|^2 &= \|P\mathbf{y} + (I - P)\mathbf{y}\|^2 \\ &= \|P\mathbf{y}\|^2 + \|(I - P)\mathbf{y}\|^2 + 2\mathbf{y}^T P^T (I - P)\mathbf{y} \\ &= \|P\mathbf{y}\|^2 + \|(I - P)\mathbf{y}\|^2 + 2\mathbf{y}^T P(I - P)\mathbf{y} \\ &= \|P\mathbf{y}\|^2 + \|(I - P)\mathbf{y}\|^2 + 2\mathbf{y}^T P\mathbf{y} - 2\mathbf{y}^T P^2\mathbf{y} \\ &= \|P\mathbf{y}\|^2 + \|(I - P)\mathbf{y}\|^2 + 2\mathbf{y}^T P\mathbf{y} - 2\mathbf{y}^T P\mathbf{y} \\ &= \|P\mathbf{y}\|^2 + \|(I - P)\mathbf{y}\|^2\end{aligned}$$

The second step could also be done by saying that  $P \perp (I - P) \rightarrow P^T(I - P) = 0$ .  
Now, a norm is always positive, hence:

$$\begin{aligned}\|P\mathbf{y}\|^2 &= \|\mathbf{y}\|^2 - \|(I - P)\mathbf{y}\|^2 \\ &\leq \|\mathbf{y}\|^2 \quad \text{as we wanted to show}\end{aligned}$$

Hence, the contraction property is proved.

ii)

The preservation of the pair-wise distances as much as possible is given by:

$$\min \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{x}_i - \mathbf{x}_j\|^2 - \|P\mathbf{x}_i - P\mathbf{x}_j\|^2$$

Expanding the norms, we have:

$$\begin{aligned}
\min \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{x}_i - \mathbf{x}_j\|^2 - \|P\mathbf{x}_i - P\mathbf{x}_j\|^2 &= \min \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j) - (P\mathbf{x}_i - P\mathbf{x}_j)^T (P\mathbf{x}_i - P\mathbf{x}_j) \\
&= \min \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j) - (P(\mathbf{x}_i - \mathbf{x}_j))^T (P(\mathbf{x}_i - \mathbf{x}_j)) \\
&= \min \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j) - (\mathbf{x}_i - \mathbf{x}_j)^T P^T P (\mathbf{x}_i - \mathbf{x}_j) \\
&= \min \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j) - (\mathbf{x}_i - \mathbf{x}_j)^T P (\mathbf{x}_i - \mathbf{x}_j) \\
&= \min \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i - \mathbf{x}_j)^T ((\mathbf{x}_i - \mathbf{x}_j) - P(\mathbf{x}_i - \mathbf{x}_j)) \\
&= \min \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i - \mathbf{x}_j)^T (I - P) (\mathbf{x}_i - \mathbf{x}_j) \\
&= \min \sum_{i,j=1}^n (\mathbf{x}_i - \mathbf{x}_j)^T (I - P) (\mathbf{x}_i - \mathbf{x}_j)
\end{aligned}$$

As we wanted to show.

**Note:** Again, we used the idempotent and symmetric properties of the projection matrix  $P$ .

iii)

First, let's consider the scatter matrix of pairwise differences  $\mathbf{S}_3$ :

$$\mathbf{S}_3 = \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i - \mathbf{x}_j) (\mathbf{x}_i - \mathbf{x}_j)^T$$

Contrary to the previous scatter matrices, this one considers the pairwise differences between all data points, and not the differences between the data points and the mean.

Now, using the cyclic and scalar properties of the trace, we have that:

$$\begin{aligned}
\min \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i - \mathbf{x}_j)^T (I - P) (\mathbf{x}_i - \mathbf{x}_j) &= \min \sum_{i=1}^n \sum_{j=1}^n \text{Tr} \left( (\mathbf{x}_i - \mathbf{x}_j)^T (I - P) (\mathbf{x}_i - \mathbf{x}_j) \right) \\
&= \min \text{Tr} \left( \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i - \mathbf{x}_j) (\mathbf{x}_i - \mathbf{x}_j)^T (I - P) \right) \\
&= \min \text{Tr} (\mathbf{S}_3 (I - P)) \quad \text{because } \mathbf{S}_3 = \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i - \mathbf{x}_j) (\mathbf{x}_i - \mathbf{x}_j)^T
\end{aligned}$$

Now, we need to introduce a scaling factor due to the double sum in the scatter matrix  $\mathbf{S}_3$ , because we are summing over all pairs  $(i, j)$ , and we are counting each pair twice, one for  $(i, j)$  and another for



$(j, i)$ . Also, we need to account for the scaling with the number of data points, since each individual datapoints is paired with all the others. Therefore, the scaling factor is given by  $2n$ , leading to:

$$\min Tr(\mathbf{S}_3(\mathbf{I} - \mathbf{P})) \rightarrow \min 2n Tr(\mathbf{S}_3(\mathbf{I} - \mathbf{P}))$$

As we wanted to show.

**d)**

The three formulations presented in the previous questions can be reduced, as shown before, to:

$$\begin{aligned} \max \operatorname{Tr}(\mathbf{S}_1 \mathbf{P}), \\ \min \operatorname{Tr}(\mathbf{S}_2(\mathbf{I} - \mathbf{P})), \\ \min \operatorname{Tr}(\mathbf{S}_3(\mathbf{I} - \mathbf{P})), \end{aligned}$$

where  $\mathbf{P}$  is the unknown projection matrix.

- The first formulation maximizes the scatter of the data projected onto the subspace.
- The second formulation minimizes the reconstruction error of the data.
- The third formulation minimizes the loss in pairwise distances of the data.

Since the trace operator returns a scalar and the identity matrix is symmetric, we can rewrite the second and third formulations as:

$$\begin{aligned} \max \operatorname{Tr}(\mathbf{S}_2(\mathbf{I} - \mathbf{P})), \\ \max \operatorname{Tr}(\mathbf{S}_3(\mathbf{I} - \mathbf{P})). \end{aligned}$$

Now, since we are maximizing the trace of a matrix, constant terms can be ignored. Specifically, in the second and third formulations,  $\operatorname{Tr}(-\mathbf{S}_2 \mathbf{I})$  and  $\operatorname{Tr}(-\mathbf{S}_3 \mathbf{I})$  are constant terms that do not explicitly affect the maximization with respect to  $\mathbf{P}$ . Hence, we can equivalently rewrite them as equivalent formulations:

$$\begin{aligned} \max \operatorname{Tr}(\mathbf{S}_1 \mathbf{P}), \\ \max \operatorname{Tr}(\mathbf{S}_2 \mathbf{P}), \\ \max \operatorname{Tr}(\mathbf{S}_3 \mathbf{P}). \end{aligned}$$

This shows that all three formulations reduce to maximizing the trace of a product of the projection matrix  $\mathbf{P}$  and a scatter matrix  $\mathbf{S}$  (which could be  $\mathbf{S}_1$ ,  $\mathbf{S}_2$ , or  $\mathbf{S}_3$ ). The solution to all three formulations is therefore the same: the matrix  $\mathbf{P}$  that maximizes the trace is the one that projects the data onto the subspace spanned by the top  $k$  eigenvectors of the covariance matrix (or scatter matrix).

Thus, the three formulations are equivalent, and they all lead to the same solution: the principal components found via PCA.

e)

The covariance matrix of the data is given by:

$$\mathbf{S} = \mathbf{V}\mathbf{D}\mathbf{V}^T$$

Where  $\mathbf{V}$  is the matrix of eigenvectors and  $\mathbf{D}$  is the diagonal matrix of eigenvalues. Those eigenvectors have the direction of the principal components, that is, the directions of the maximum variance of the data. Hence, the dotted lines in the plot represent the directions of the principal components and, thus, the directions of the eigenvectors of  $\mathbf{V}$ , which are the columns of  $\mathbf{V}$ . In the plot, the yellow line corresponds to the direction of  $V_1$  and the brown line corresponds to the direction of  $V_2$ .

Now, for any given space that is spanned by the column space of  $\mathbf{A}$ , the projection matrix  $\mathbf{P}$  is given by:

$$\mathbf{P} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$$

Since, in the case of PCA, the matrix  $\mathbf{A}$  is the matrix of eigenvectors, the projection matrix  $\mathbf{P}$  is given by:

$$\mathbf{P} = \mathbf{V}_k \mathbf{V}_k^T$$

Where  $\mathbf{V}_k$  is the matrix of the first  $k$  eigenvectors of the covariance matrix. This is because, the eigenvectors are orthogonal to each other, hence:

$$\mathbf{P} = \mathbf{V}_k (\mathbf{V}_k^T \mathbf{V}_k)^{-1} \mathbf{V}_k^T = \mathbf{V}_k (\mathbf{I})^{-1} \mathbf{V}_k^T = \mathbf{V}_k \mathbf{V}_k^T$$

Assuming  $k = 1$ , the projection matrix becomes:

$$\mathbf{P} = \mathbf{V}_1 \mathbf{V}_1^T$$

Then, the projection of the data onto the first principal component by adding the bias back is given by:

$$\mathbf{P}\mathbf{x}_i = \mathbf{V}_1 \mathbf{V}_1^T \mathbf{x}_i + \bar{\mathbf{x}}$$

Not considering the bias, the projection of the data onto the first principal component is given by:

$$\mathbf{P}\mathbf{x}_i = \mathbf{V}_1 \mathbf{V}_1^T \mathbf{x}_i$$

If we consider the correction of the bias, the projection will be a parallel line to the yellow dotted line, in which the bias is the distance between the origin and the mean of the data. If we do not consider the

correction of the bias, the points will be projected onto the yellow dotted line.

Plotting the projection of the data onto the first principal component, we have:

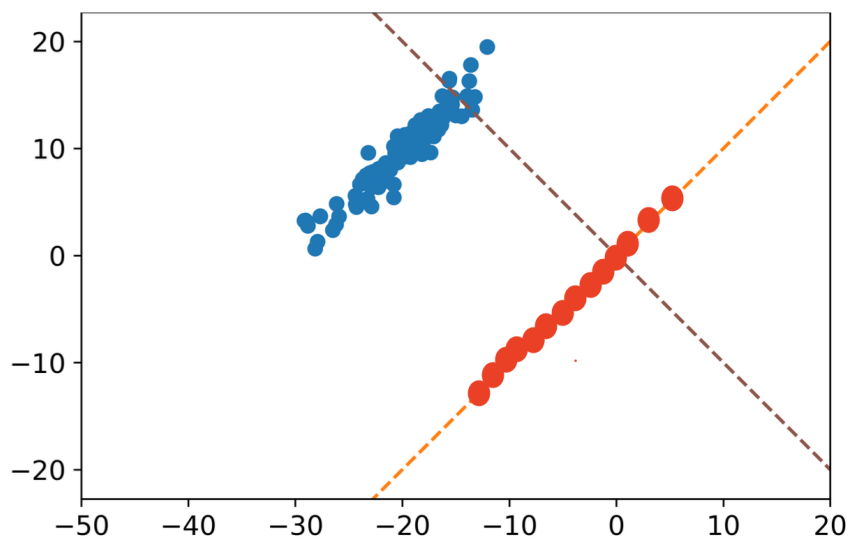


Fig. 1: Projection of the data onto the first principal component

## 2 Probabilistic PCA - A general latent space distribution

a)

Since

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$$

Then we can say that:

$$\epsilon_x \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$$

Now, since  $\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \epsilon_x$  and  $\mathbf{z} = \mathbf{m} + \epsilon_z$ , we can rewrite the expression as:

$$\begin{aligned}\mathbf{x} &= \mathbf{W}(\mathbf{m} + \epsilon_z) + \boldsymbol{\mu} + \epsilon_x \\ &= \mathbf{W}\mathbf{m} + \mathbf{W}\epsilon_z + \boldsymbol{\mu} + \epsilon_x\end{aligned}$$

Now, since  $\mathbf{W}\epsilon_z$  is a linear transformation of a Gaussian random variable, it is also Gaussian. Therefore,  $\mathbf{W}\epsilon_z \sim \mathcal{N}(0, \mathbf{W}\Sigma\mathbf{W}^T)$  (this will be shown in the next question).

With this, since the sum of two Gaussian random variables is also Gaussian, we have that:

$$\mathbf{x} \sim \mathcal{N}(\mathbf{W}\mathbf{m} + \boldsymbol{\mu}, \mathbf{W}\Sigma\mathbf{W}^T + \sigma^2 \mathbf{I})$$

Note: We can ignore the  $\mathbf{W}(\mathbf{m} + \boldsymbol{\mu})$  because they are constants and do not affect the type of distribution of  $\mathbf{x}$ .

b)

To find the expectation of the variable  $\mathbf{x}$  by taking into account the linearity of the expectation operator, we have that:

$$\begin{aligned}E[\mathbf{x}] &= E[\mathbf{W}\mathbf{m} + \boldsymbol{\mu} + \epsilon_x + \mathbf{W}\epsilon_z] \\ &= E[\mathbf{W}\mathbf{m}] + E[\boldsymbol{\mu}] + E[\epsilon_x] + E[\mathbf{W}\epsilon_z] \\ &= \mathbf{W}\mathbf{m} + \boldsymbol{\mu} + E[\epsilon_x] + \mathbf{W}E[\epsilon_z] \\ &= \mathbf{W}\mathbf{m} + \boldsymbol{\mu}\end{aligned}$$

This, because the expectancy of a constant is the constant itself in the case of  $\mathbf{W}\mathbf{m}$  and  $\boldsymbol{\mu}$ , and the

expectancy of a Gaussian random variable is the mean of the Gaussian distribution, which is 0 in the case of  $\epsilon_x$ , since  $\epsilon_x \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$  and in the case of  $\mathbf{W}\epsilon_z$ , since  $\epsilon_z \sim \mathcal{N}(0, \Sigma)$ .

**c)**

To find the covariance of the variable  $x$ , we have that:

$$\begin{aligned}
Cov[\mathbf{x}, \mathbf{x}] &= Var[\mathbf{x}] \\
&= Var[\mathbf{W}\mathbf{m} + \boldsymbol{\mu} + \epsilon_x + \mathbf{W}\epsilon_z] \\
&= Var[\mathbf{W}\mathbf{m}] + Var[\boldsymbol{\mu}] + Var[\epsilon_x] + Var[\mathbf{W}\epsilon_z] \\
&= \mathbf{W}Var[\mathbf{m}]\mathbf{W}^T + Var[\boldsymbol{\mu}] + Var[\epsilon_x] + \mathbf{W}Var[\epsilon_z]\mathbf{W}^T \\
&= \mathbf{W}Var[\mathbf{m}]\mathbf{W}^T + Var[\epsilon_x] + \mathbf{W}Var[\epsilon_z]\mathbf{W}^T \\
&= \mathbf{W}Var[\mathbf{m}]\mathbf{W}^T + \sigma^2 \mathbf{I} + \mathbf{W}\Sigma\mathbf{W}^T \\
&= \sigma^2 \mathbf{I} + \mathbf{W}\Sigma\mathbf{W}^T \quad \text{since } Var[\mathbf{m}] = Var[\boldsymbol{\mu}] = 0 \text{ because they are constants} \\
&= \sigma^2 \mathbf{I} + \mathbf{W}\Sigma\mathbf{W}^T \quad \text{since } Var[\epsilon_x] = \sigma^2 \mathbf{I} \text{ and } Var[\epsilon_z] = \Sigma
\end{aligned}$$

This, because the covariance of a Gaussian random variable is the covariance matrix of the Gaussian distribution.

**d)**

To match the previous expression for the General Gaussian prior:

$$\mathbf{x} \sim \mathcal{N}(\mathbf{W}\mathbf{m} + \boldsymbol{\mu}, \mathbf{W}\Sigma\mathbf{W}^T + \sigma^2 \mathbf{I})$$

In the form:

$$\mathbf{x} \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\mathbf{W}}\tilde{\mathbf{W}}^T + \sigma^2 \mathbf{I})$$

We have that:

$$\begin{aligned}
\tilde{\boldsymbol{\mu}} &= \mathbf{W}\mathbf{m} + \boldsymbol{\mu} \\
\tilde{\mathbf{W}} &= \mathbf{W}\Sigma^{1/2}
\end{aligned}$$

The  $\tilde{\boldsymbol{\mu}}$  is easy to check because it is the mean of the Gaussian distribution, and the  $\tilde{\mathbf{W}}$  is the square root of the covariance matrix of the Gaussian distribution, since:

$$\begin{aligned}
\mathbf{W}\Sigma\mathbf{W}^T + \sigma^2\mathbf{I} &= \mathbf{W}\Sigma^{1/2}\Sigma^{1/2}\mathbf{W}^T + \sigma^2\mathbf{I} \\
&= \mathbf{W}\Sigma^{1/2}\Sigma^{1/2T}\mathbf{W}^T + \sigma^2\mathbf{I} \quad \text{since } \Sigma \text{ is symmetric} \\
&= (\mathbf{W}\Sigma^{1/2})(\mathbf{W}\Sigma^{1/2})^T + \sigma^2\mathbf{I}
\end{aligned}$$

### 3 Mixtures of Experts

a)

Considering that  $z_n$  is one-hot encoded, then we assign 1 to the index corresponding to the expert that is responsible for the data point  $\mathbf{x}_n$ , i.e., the one with the highest probability, and 0 to the other indices. In resume, we need to find the index  $k$  that maximizes the probability  $p(z_n = k|\mathbf{x}_n, \Phi) = \pi_{nk}$ .

This gives that:

$$z_n = \underset{j}{\operatorname{argmax}} \pi_{nj} = \begin{cases} 1, & \text{if } k = \underset{j}{\operatorname{argmax}} \pi_{nj} = \underset{j}{\operatorname{argmax}} \frac{\exp(\phi_j^T \mathbf{x}_n)}{\sum_{l=1}^K \exp(\phi_l^T \mathbf{x}_n)} \\ 0, & \text{otherwise} \end{cases}$$

b)

The likelihood of the data is given by:

$$\begin{aligned}
p(\mathbf{y}|\mathbf{X}, \Theta, \Phi) &= \prod_{n=1}^N p(\mathbf{y}_n|\mathbf{x}_n, \Theta, \Phi) \quad \text{assuming i.i.d.} \\
&= \prod_{n=1}^N \sum_{k=1}^K p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \Theta) p(z_n = k|\mathbf{x}_n, \Phi) \quad \text{marginalizing } z_n \text{ over experts} \\
&= \prod_{n=1}^N \sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \Theta)
\end{aligned}$$

Which, if we expand the terms, gives:

$$\begin{aligned}
p(\mathbf{y}|\mathbf{X}, \Theta, \Phi) &= \prod_{n=1}^N \sum_{k=1}^K \frac{\exp(\phi_k^T \mathbf{x}_n)}{\sum_{l=1}^K \exp(\phi_l^T \mathbf{x}_n)} \exp(\theta_k^T \mathbf{x}_n) \exp(-\exp(\theta_k^T \mathbf{x}_n) \mathbf{y}_n) \\
&= \prod_{n=1}^N \sum_{k=1}^K \frac{\exp(\phi_k^T \mathbf{x}_n)}{\sum_{l=1}^K \exp(\phi_l^T \mathbf{x}_n)} \exp(\theta_k^T \mathbf{x}_n - \exp(\theta_k^T \mathbf{x}_n) \mathbf{y}_n)
\end{aligned}$$

The log-likelihood, without the expanding terms, is given by:

$$\log p(\mathbf{y}|\mathbf{X}, \Theta, \Phi) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | z_n = k, \mathbf{x}_n, \Theta)$$

Expanding the terms, we have:

$$\log p(\mathbf{y}|\mathbf{X}, \Theta, \Phi) = \sum_{n=1}^N \log \sum_{k=1}^K \frac{\exp(\phi_k^T \mathbf{x}_n)}{\sum_{l=1}^K \exp(\phi_l^T \mathbf{x}_n)} \exp(\theta_k^T \mathbf{x}_n - \exp(\theta_k^T \mathbf{x}_n) \mathbf{y}_n)$$

c)

The responsibility of the expert  $i$  for the data point  $\mathbf{x}_n$  is given by:

$$\begin{aligned} r_{ni} &= p(z_{ni} = 1 | \mathbf{x}_n, \Theta, \Phi) \\ &= \frac{p(\mathbf{y}_n | \mathbf{x}_n, z_{ni} = 1, \Theta) p(z_{ni} = 1 | \mathbf{x}_n, \Phi)}{p(\mathbf{y}_n | \mathbf{x}_n, \Theta, \Phi)} \\ &= \frac{p(\mathbf{y}_n | z_{ni} = 1, \mathbf{x}_n, \Theta) p(z_{ni} = 1 | \mathbf{x}_n, \Phi)}{\sum_{j=1}^K p(\mathbf{y}_n | z_{nj} = 1, \mathbf{x}_n, \Theta) p(z_{nj} = 1 | \mathbf{x}_n, \Phi)} \\ &= \frac{\pi_{ni} p(\mathbf{y}_n | z_{ni} = 1, \mathbf{x}_n, \Theta)}{\sum_{j=1}^K \pi_{nj} p(\mathbf{y}_n | z_{nj} = 1, \mathbf{x}_n, \Theta)} \end{aligned}$$

By expanding the terms, we have:

$$r_{ni} = \frac{\frac{\exp(\phi_i^T \mathbf{x}_n)}{\sum_{l=1}^K \exp(\phi_l^T \mathbf{x}_n)} \exp(\theta_i^T \mathbf{x}_n - \exp(\theta_i^T \mathbf{x}_n) \mathbf{y}_n)}{\sum_{j=1}^K \frac{\exp(\phi_j^T \mathbf{x}_n)}{\sum_{l=1}^K \exp(\phi_l^T \mathbf{x}_n)} \exp(\theta_j^T \mathbf{x}_n - \exp(\theta_j^T \mathbf{x}_n) \mathbf{y}_n)}$$

Which can be simplified to:

$$r_{ni} = \frac{\exp(\phi_i^T \mathbf{x}_n) \exp(\theta_i^T \mathbf{x}_n - \exp(\theta_i^T \mathbf{x}_n) \mathbf{y}_n)}{\sum_{j=1}^K \exp(\phi_j^T \mathbf{x}_n) \exp(\theta_j^T \mathbf{x}_n - \exp(\theta_j^T \mathbf{x}_n) \mathbf{y}_n)}$$

d)

As shown before, the likelihood is given by:



$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\Theta}, \Phi) = \prod_{n=1}^N \sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\Theta})$$

By taking the hint into account to derivate the log-likelihood, we have that:

$$\begin{aligned} \frac{\partial \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\Theta}, \Phi)}{\partial \boldsymbol{\theta}_i} &= \frac{\partial}{\partial \boldsymbol{\theta}_i} \sum_{n=1}^N \log \sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k) \\ &= \sum_{n=1}^N \frac{\partial}{\partial \boldsymbol{\theta}_i} \log \sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k) \\ &= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k)} \frac{\partial}{\partial \boldsymbol{\theta}_i} \sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k) \\ &= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k)} \sum_{k=1}^K \frac{\partial}{\partial \boldsymbol{\theta}_i} p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k) \end{aligned}$$

Since,  $p(z_n = k|\mathbf{x}_n, \Phi)$  does not depend on  $\boldsymbol{\theta}_i$ , we have that:

$$\frac{\partial \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\Theta}, \Phi)}{\partial \boldsymbol{\theta}_i} = \sum_{n=1}^N \frac{1}{\sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k)} \sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) \frac{\partial}{\partial \boldsymbol{\theta}_i} p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k)$$

Now, taking the derivative in respect of  $\phi_i$ :

$$\begin{aligned} \frac{\partial \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\Theta}, \Phi)}{\partial \phi_i} &= \frac{\partial}{\partial \phi_i} \sum_{n=1}^N \log \sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k) \\ &= \sum_{n=1}^N \frac{\partial}{\partial \phi_i} \log \sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k) \\ &= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k)} \frac{\partial}{\partial \phi_i} \sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k) \\ &= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k)} \sum_{k=1}^K \frac{\partial}{\partial \phi_i} p(z_n = k|\mathbf{x}_n, \Phi) p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k) \end{aligned}$$

Now, since  $p(\mathbf{y}_n|z_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k)$  does not depend on  $\phi_i$ , we have that:

$$\frac{\partial \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\Theta}, \Phi)}{\partial \phi_i} = \sum_{n=1}^N \frac{1}{\sum_{k=1}^K p(\mathbf{z}_n = k|\mathbf{x}_n, \Phi)p(\mathbf{y}_n|\mathbf{z}_n = k, \mathbf{x}_n, \Theta)} \sum_{k=1}^K p(\mathbf{y}_n|\mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\theta}_k) \frac{\partial}{\partial \phi_i} p(\mathbf{z}_n = k|\mathbf{x}_n, \Phi)$$

**Note:** The responsibilities only appear when taking the explicitly derivation of the probabilities, since it will depend on the values of the probabilities. This is shown in the next question.

e)

Now, inserting the explicit expressions for the probabilities, we have that:

$$\begin{aligned}\frac{\partial \log p(\mathbf{y}|\mathbf{X}, \Theta, \Phi)}{\partial \theta_i} &= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta)} \frac{\partial}{\partial \theta_i} \sum_{k=1}^K p(\mathbf{z}_n = k | \mathbf{x}_n, \Phi) p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \theta_k) \\ &= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta)} \frac{\partial}{\partial \theta_i} \sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta)\end{aligned}$$

Since  $\pi_{nk}$  does not depend on  $\theta_i$ , we have that:

$$\begin{aligned}\frac{\partial \log p(\mathbf{y}|\mathbf{X}, \Theta, \Phi)}{\partial \theta_i} &= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta)} \sum_{k=1}^K \pi_{nk} \frac{\partial}{\partial \theta_i} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta) \\ &= \sum_{n=1}^N \frac{\pi_{nk}}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta)} \sum_{k=1}^K \frac{\partial}{\partial \theta_i} (\lambda \exp(-\lambda \mathbf{y}_n)) \\ &= \sum_{n=1}^N \frac{\pi_{nk}}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta)} \sum_{k=1}^K \frac{\partial}{\partial \lambda} (\lambda \exp(-\lambda \mathbf{y}_n)) \frac{\partial \lambda}{\partial \theta_i} \quad \text{by the chain rule} \\ &= \sum_{n=1}^N \frac{\pi_{nk}}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta)} \sum_{k=1}^K (\exp(-\lambda \mathbf{y}_n) - \lambda \mathbf{y}_n \exp(-\lambda \mathbf{y}_n)) \frac{\partial \lambda}{\partial \theta_i} \\ &= \sum_{n=1}^N \frac{\pi_{nk}}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta)} \sum_{k=1}^K (\exp(-\lambda \mathbf{y}_n) - \lambda \mathbf{y}_n \exp(-\lambda \mathbf{y}_n)) \frac{\partial}{\partial \theta_i} \exp(\theta_k^T \mathbf{x}_n) \\ &= \sum_{n=1}^N \frac{\pi_{nk}}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta)} \sum_{k=1}^K (\exp(-\lambda \mathbf{y}_n) - \lambda \mathbf{y}_n \exp(-\lambda \mathbf{y}_n)) \exp(\theta_k^T \mathbf{x}_n) \mathbf{x}_n^T \delta_{ik} \\ &= \sum_{n=1}^N \frac{\pi_{ni}}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta)} (\exp(-\lambda \mathbf{y}_n) - \lambda \mathbf{y}_n \exp(-\lambda \mathbf{y}_n)) \exp(\theta_i^T \mathbf{x}_n) \mathbf{x}_n^T \\ &= \sum_{n=1}^N \frac{\pi_{ni}}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta)} \lambda \exp(-\lambda \mathbf{y}_n) (\mathbb{1} - \lambda \mathbf{y}_n) \mathbf{x}_n^T \\ &= \sum_{n=1}^N \frac{\pi_{ni} p(\mathbf{y}_n | \mathbf{z}_n = i, \mathbf{x}_n, \theta_i)}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \Theta)} (\mathbb{1} - \lambda \mathbf{y}_n) \mathbf{x}_n^T \quad \text{since } p(\mathbf{y}_n | \mathbf{z}_n = i, \mathbf{x}_n, \theta_i) = \lambda \exp(-\lambda \mathbf{y}_n) \\ &= \sum_{n=1}^N r_{ni} \mathbf{x}_n^T (\mathbb{1} - \lambda \mathbf{y}_n)\end{aligned}$$

Now, for  $\phi_i$ :

$$\begin{aligned}
\frac{\partial \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\Theta}, \boldsymbol{\Phi})}{\partial \phi_i} &= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Phi})} \frac{\partial}{\partial \phi_i} \sum_{k=1}^K p(\mathbf{z}_n = k | \mathbf{x}_n, \boldsymbol{\Phi}) p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Phi}) \\
&= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Theta})} \frac{\partial}{\partial \phi_i} \sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Theta}) \\
&= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Theta})} \sum_{k=1}^K \frac{\partial}{\partial \phi_i} \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Theta}) \\
&= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Theta})} \sum_{k=1}^K p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Theta}) \frac{\partial}{\partial \phi_i} \pi_{nk} \\
&= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Theta})} \sum_{k=1}^K p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Theta}) \frac{\partial}{\partial \phi_i} \frac{\exp(\boldsymbol{\phi}_k^T \mathbf{x}_n)}{\sum_{l=1}^K \exp(\boldsymbol{\phi}_l^T \mathbf{x}_n)} \\
&= \sum_{n=1}^N \frac{p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Theta}) \pi_{ni} (1 - \pi_{ni}) \mathbf{x}_n^T - \sum_{k \neq i}^K p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Theta}) \pi_{nk} \pi_{ni} \mathbf{x}_n^T}{\sum_{k=1}^K \pi_{nk} p(\mathbf{y}_n | \mathbf{z}_n = k, \mathbf{x}_n, \boldsymbol{\Theta})} \\
&= \sum_{n=1}^N \left[ r_{ni} (1 - \pi_{ni}) \mathbf{x}_n^T - \sum_{k \neq i}^K r_{nk} \pi_{ni} \mathbf{x}_n^T \right] \\
&= \sum_{n=1}^N \left[ r_{ni} (1 - \pi_{ni}) \mathbf{x}_n^T - \pi_{ni} \mathbf{x}_n^T \sum_{k \neq i}^K r_{nk} \right] \\
&= \sum_{n=1}^N [r_{ni} (1 - \pi_{ni}) \mathbf{x}_n^T - \pi_{ni} \mathbf{x}_n^T (1 - r_{ni})] \quad \text{since } \sum_{k=1}^K r_{nk} = 1 \rightarrow \sum_{k \neq i}^K r_{nk} = 1 - r_{ni} \\
&= \sum_{n=1}^N [r_{ni} - \pi_{ni}] \mathbf{x}_n^T
\end{aligned}$$

f)

Writing down an iterative algorithm that maximizes the log-probability by jointly optimizing the parameters  $\Theta$  and  $\Phi$ :

---

**Algorithm 1** Joint optimization of  $\Theta$  and  $\Phi$  in Mixtures of Experts

---

**Input:** Data  $X$ , labels  $Y$ , parameters  $\Theta$ ,  $\Phi$ , learning rate  $\alpha$ , tolerance  $\epsilon$ , maximum iterations

**Output:** Updated parameters  $\Theta$ ,  $\Phi$

Initialize parameters  $\Theta$  and  $\Phi$

**while** not converged **do**

**E-step:** Compute responsibilities  $r_{nk}$ :

$$r_{nk} = \frac{p(\mathbf{y}_n | \mathbf{x}_n, \boldsymbol{\theta}_k) \pi_{nk}}{\sum_{j=1}^K p(\mathbf{y}_n | \mathbf{x}_n, \boldsymbol{\theta}_j) \pi_{nj}}$$

where  $\pi_{nk} = p(\mathbf{z}_n = k | \mathbf{x}_n, \Phi) = \frac{\exp(\phi_k^T \mathbf{x}_n)}{\sum_{j=1}^K \exp(\phi_j^T \mathbf{x}_n)}$

**M-step:** Update  $\Theta$  and  $\Phi$  using the gradients:

$$\theta_i^T \leftarrow \theta_i^T + \alpha \sum_{n=1}^N r_{ni} \mathbf{x}_n^T (\mathbb{1} - \lambda \mathbf{y}_n)$$

$$\phi_i^T \leftarrow \phi_i^T + \alpha \sum_{n=1}^N [r_{ni} - \pi_{ni}] \mathbf{x}_n^T$$

**Check for convergence:**

    Compute the change in log-likelihood:

$$\Delta \mathcal{L} = |\mathcal{L}^{(t+1)} - \mathcal{L}^{(t)}|$$

**if**  $\Delta \mathcal{L} < \epsilon$  **or** maximum iterations reached **then**

**Terminate:** The algorithm has converged when the log-likelihood improvement is smaller than a predefined threshold  $\epsilon$ , indicating that further updates provide negligible improvement.

**Return** optimized parameters  $\Theta$  and  $\Phi$ .

**else**

        Continue to the next iteration by updating  $t \leftarrow t + 1$ .

**end if**

**end while**

---

**g)**

If instead of having a single expert, we have multiple experts, then the final prediction  $\hat{\mathbf{y}}$  is given by the weighted sum of the predictions of each expert, where the weights are given by the responsibilities  $r_{nk}$ :

$$\hat{\mathbf{y}} = \sum_{k=1}^K r_{nk} \hat{\mathbf{y}}_{nk}$$

This  $\hat{\mathbf{y}}_{nk}$  is the prediction of the expert  $k$  for the data point  $\mathbf{x}_n$ , which can be sampled from the distribution of expert  $k$ .