

Assignment Number: 2
Student Name: Gurpreet Singh
Roll Number: 150259
Date: November 8, 2017

Part 1

We know that Gaussian-Gaussian has the conjugacy property, and hence, we can say that the marginal probability of \mathbf{x} i.e. $\mathbb{P}[\mathbf{x}]$

Therefore, we can write the marginal probability of \mathbf{x} as a gaussian distribution with mean $\boldsymbol{\mu}'$ and covariance matrix Λ .

$$\begin{aligned}\mathbb{P}[\mathbf{x} | \boldsymbol{\mu}, W, \sigma] &= \int_{\mathbf{z}} \mathbb{P}[\mathbf{x} | \mathbf{z}, \boldsymbol{\mu}, W, \sigma] \mathbb{P}[\mathbf{z} | 0, \mathbf{I}] d\mathbf{z} \\ &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}', \Lambda)\end{aligned}$$

Also, we know that $\mathbf{x} = W\mathbf{z} + \boldsymbol{\mu} + \epsilon$, and from the properties of Gaussian, we know $\boldsymbol{\mu}' = \mathbb{E}[\mathbf{x}]$ and $\Lambda = \text{cov}(\mathbf{x})$. Hence

$$\begin{aligned}\boldsymbol{\mu}' &= \mathbb{E}[\mathbf{x}] \\ &= \mathbb{E}[W\mathbf{z} + \boldsymbol{\mu} + \epsilon] \\ &= W(\mathbb{E}[\mathbf{z}] = 0) + (\mathbb{E}[\boldsymbol{\mu}] = \boldsymbol{\mu}) + (\mathbb{E}[\epsilon] = 0) \\ &= \boldsymbol{\mu}\end{aligned}$$

$$\begin{aligned}\Lambda' &= \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}]) (\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] \\ &= \mathbb{E}[(W\mathbf{z} + \epsilon) (W\mathbf{z} + \epsilon)^T] \\ &= \mathbb{E}[W\mathbf{z}\mathbf{z}^T W^T + W\mathbf{z}\epsilon^T + \epsilon W\mathbf{z}^T + \epsilon\epsilon^T] \\ &= W\mathbb{E}[\mathbf{z}\mathbf{z}^T] W^T + 0 + 0 + \mathbb{E}[\epsilon\epsilon^T] \\ &= W\mathbf{I}\mathbf{I}^T W^T + \sigma^2 \mathbf{I} \\ &= WW^T + \sigma^2 \mathbf{I}\end{aligned}$$

Hence, we can write $\mathbb{P}[\mathbf{x}] = \mathcal{N}(\boldsymbol{\mu}, C)$ where $C = WW^T + \sigma^2 \mathbf{I}$

Part 2

Note: I will use $C = WW^T + \sigma^2 \mathbf{I}$ for the rest of the question

Since all the sample points are independent, we can write $\mathbb{P}[X] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i]$. Therefore

$$\begin{aligned}
\mathbb{P}[X \mid \boldsymbol{\mu}, W, \sigma] &= \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i \mid \boldsymbol{\mu}, W, \sigma] \\
&= \prod_{i=1}^n \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}, C) \\
\Rightarrow \mathbb{P}[X \mid \boldsymbol{\mu}, W, \sigma] &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi \|C\|}} e^{\frac{-1}{2}(\mathbf{x}^i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}^i - \boldsymbol{\mu})}
\end{aligned}$$

Part 3

We can write the MLE estimate of $\boldsymbol{\mu}$ as

$$\begin{aligned}
\boldsymbol{\mu}^{MLE} &= \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^d} \mathbb{P}[X \mid \boldsymbol{\mu}, W, \sigma] \\
&= \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^d} \log(\mathbb{P}[X \mid \boldsymbol{\mu}, W, \sigma]) \\
&= \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^d} \log\left(\prod_{i=1}^n \mathbb{P}[\mathbf{x}^i \mid \boldsymbol{\mu}, W, \sigma]\right) \\
&= \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^d} \sum_{i=1}^n \log(\mathbb{P}[\mathbf{x}^i \mid \boldsymbol{\mu}, W, \sigma]) \\
&= \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} \sum_{i=1}^n \frac{1}{2} (\mathbf{x}^i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}^i - \boldsymbol{\mu}) + \frac{1}{2} \log(2\pi \|C\|) \\
&= \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} \sum_{i=1}^n \frac{1}{2} (\mathbf{x}^i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}^i - \boldsymbol{\mu})
\end{aligned}$$

We can simply differentiate the RHS term, and using differentials of matrices, we get

$$\begin{aligned}
\frac{\delta(RHS)}{\delta(\boldsymbol{\mu})} &= \frac{\delta\left(\sum_{i=1}^n \frac{1}{2} (\mathbf{x}^i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}^i - \boldsymbol{\mu})\right)}{\delta(\boldsymbol{\mu})} = 0 \\
\Rightarrow \sum_{i=1}^n C^{-1} (\mathbf{x}^i - \boldsymbol{\mu}) &= 0 \\
\Rightarrow \boldsymbol{\mu}^{MLE} &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}^i
\end{aligned}$$

Therefore, the MLE estimate of $\boldsymbol{\mu}$ is, as expected, the empirical mean of all the data points. Hence we use this to centralize for the PCA algorithm.