

I worked on it alone, Group E.

1 Problem 1

1. In order to obtain the update rules for the asked parameters, we derive the expected value of the complete data log-likelihood with regard to the respective parameters and set it zero. In the same time, the term $\gamma(z_{nk})$ is kept fixed and treated as a constant. In order to solve the entire problem, we have to draw on Lagrangian Multipliers with the following form:

$$\begin{aligned} F &= \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{ \ln \pi_k + \ln N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \} + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left\{ \ln \pi_k - \frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right\} + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) + \text{const.}, \end{aligned}$$

where the term on the right is the introduced constraint which is necessary for Lagrangian optimization with constraints. It says: $\sum_{k=1}^K \pi_k = 1$.

Now, we can start deriving this full Lagrangian with respect to the individual parameters to end up with their update rules:

$$\begin{aligned} \frac{\partial F}{\partial \pi_k} &= \sum_{n=1}^N \frac{\gamma(z_{nk})}{\pi_k} + \lambda = 0 \\ N_k &= -\lambda \pi_k \\ \sum_{k=1}^K N_k &= -\sum_{k=1}^K \lambda \pi_k \\ N &= -\lambda \\ \pi_k &= \frac{N_k}{N}, \end{aligned}$$

with $N_k = \sum_{n=1}^N \gamma(z_{nk})$ being the number of data points associated with component k.

For μ : At first, we can simplify the F term for the terms being affected by the derivative regarding μ and also multiply the quadratic term:

$$F_\mu = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left[-\frac{1}{2} (\mathbf{x}_n^T \boldsymbol{\Sigma}_k^{-1} \mathbf{x}_n - 2 \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}_k^{-1} \mathbf{x}_n + \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k) \right] + \text{const},$$

which is then again derived for the respective parameter as such:

$$\begin{aligned}
\frac{\partial F_{\mu}}{\partial \mu_k} &= -\frac{1}{2} \sum_{n=1}^N \gamma(z_{nk}) (-2\Sigma_k^{-1} \mathbf{x}_n + 2\Sigma_k^{-1} \mu_k) = 0 \\
\sum_{n=1}^N \gamma(z_{nk}) \Sigma_k^{-1} \mathbf{x}_n &= \sum_{n=1}^N \gamma(z_{nk}) \Sigma_k^{-1} \mu_k \\
\sum_{n=1}^N \gamma(z_{nk}) \Sigma_k^{-1} \mathbf{x}_n &= N_k \Sigma_k^{-1} \mu_k \\
\mu_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n
\end{aligned}$$

Similarly for the Covariance matrix Σ_k a form of the Lagrangian with only relevant terms remaining:

$$F_{\Sigma_k} = -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left[\ln(|\Sigma_k|) + (\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) \right] + \text{const}$$

Then the derivation with setting zero arrives as:

$$\begin{aligned}
\frac{\partial F_{\Sigma_k}}{\partial \Sigma_k} &= -\frac{1}{2} \sum_{n=1}^N \gamma(z_{nk}) \left[\Sigma_k^{-1} - \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) (\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} \right] = 0 \\
\sum_{n=1}^N \gamma(z_{nk}) \Sigma_k^{-1} &= \sum_{n=1}^N \gamma(z_{nk}) \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) (\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} \\
\Sigma_k N_k \Sigma_k^{-1} \Sigma_k &= \sum_{n=1}^N \gamma(z_{nk}) \Sigma_k \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) (\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} \Sigma_k \\
\Sigma_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k) (\mathbf{x}_n - \mu_k)^T
\end{aligned}$$

I used the following formulas in this order: matrix cookbook equation (57): $\frac{\partial \ln |\det(\mathbf{X})|}{\partial \mathbf{X}} = (\mathbf{X}^{-1})^T = (\mathbf{X}^T)^{-1}$. Moreover, matrix cookbook equation (61): $\frac{\partial \mathbf{a}^T \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -\mathbf{X}^{-T} \mathbf{a} \mathbf{b}^T \mathbf{X}^{-T}$ and $\Sigma_k \Sigma_k^{-1} = I$. In the third row, Σ_k is multiplied from the left and the right to eliminate the inverse of this term.

- Supposing, the covariance matrices are constrained to have the same value, the update rules of μ and π do not change, as they are independent of Σ . Again, I give the respective Lagrangian with only the relevant terms remaining, now Σ instead of Σ_k :

$$F_{\Sigma} = -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left[\ln(|\Sigma|) + (\mathbf{x}_n - \mu_k)^T \Sigma^{-1} (\mathbf{x}_n - \mu_k) \right] + \text{const}$$

Analogously applying the same derivation formulas as before, except multiplying Σ with no index k from the left and the right:

$$\begin{aligned}
\frac{F_{\Sigma_k}}{\partial \Sigma} &= -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left[\Sigma^{-1} - \Sigma^{-1} (\mathbf{x}_n - \mu_k) (\mathbf{x}_n - \mu_k)^T \Sigma^{-1} \right] = 0 \\
\sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \Sigma^{-1} &= \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \Sigma^{-1} (\mathbf{x}_n - \mu_k) (\mathbf{x}_n - \mu_k)^T \Sigma^{-1} \\
\Sigma &= \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) (\mathbf{x}_n - \mu_k) (\mathbf{x}_n - \mu_k)^T,
\end{aligned}$$

with in the last row: $\sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) = \sum_{k=1}^K N_k = N$

2 Problem 2

Given the graphical model with its dependencies, the posterior distribution can be written as:

$$\begin{aligned}
 p(\boldsymbol{\theta}|\mathbf{X}) &= \frac{p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{X})} \\
 &\propto p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \\
 &= \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \\
 &\propto \ln \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta})
 \end{aligned}$$

In the E-Step of the EM Algorithm, the posterior is maximized with respect to the latent variables, while at the same time the parameters θ are kept fixed:

$$\arg \max_{\mathbf{Z}} p(\boldsymbol{\theta}|\mathbf{X}) \propto \arg \max_{\mathbf{Z}} \ln \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta})$$

As the prior is independent of the latent parameters, which are being maximized, it can be dropped. That leaves the maximization of the log-likelihood, which is the same as for Maximum Likelihood:

$$\propto \arg \max_{\mathbf{Z}} \ln \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$$

In the M-Step, the latent parameters are fixed, while at the same time maximizing the posterior with respect to the parameters θ . This time, the prior remains:

$$\arg \max_{\mathbf{Z}} p(\boldsymbol{\theta}|\mathbf{X}) \propto \arg \max_{\mathbf{Z}} \ln \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta})$$

However, as the latent parameters are not given, maximizing over the full data log-likelihood is not possible. Instead, the only derivable info on the latent parameters is with respect to the posterior distribution: $p(\mathbf{Z}|\mathbf{X}, \theta)$. Therefore, the the left term in the M-step formula can has to approximated with the expected value under the posterior distribution of the latent variables, as we can derive this information. With that, we end up as such:

$$\begin{aligned}
 \arg \max_{\theta} p(\theta|\mathbf{X}) &\propto \arg \max_{p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})} \left[\ln \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right] + \ln p(\theta) \\
 &\approx \arg \max_{\theta} \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) + \ln p(\theta)
 \end{aligned}$$

3 Problem 3

In the M-Step, update rules for π and μ are sought. At first, we formulate the log-posterior as such:

$$\begin{aligned}
 \ln p(\boldsymbol{\mu}, \boldsymbol{\pi} | \{x_n\}_{n=1}^N) &\propto \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left(\ln \boldsymbol{\pi}_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln (1 - \mu_{ki})] \right) \\
 &\quad + \sum_{k=1}^K \ln p(\boldsymbol{\mu}_k | a_k, b_k) + \ln p(\boldsymbol{\pi} | \boldsymbol{\alpha})
 \end{aligned}$$

Now, the first distribution in the second row is the Beta distribution of parameter μ_k and the second is the Dirichlet of π . Plugging those distributions in, we arrive at:

$$\begin{aligned}
&= \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left(\ln \pi_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln (1 - \mu_{ki})] \right) \\
&\quad + \sum_{k=1}^K (a_k - 1) \ln \mu_k + (b_k - 1) \ln (1 - \mu_k) + \sum_{k=1}^K (\alpha_k - 1) \ln \pi_k + \text{const},
\end{aligned}$$

whereas const entails the Beta function $B(a_k, b_k)$ term of the beta distribution and $B(a_k)$ of the Dirichlet distribution. As we derive for μ and π , they do not matter and can be summarized as a constant.

As seen in the first Problem, for Lagrange Multipliers, we need to introduce a constraint. Again, with $\sum_{k=1}^K \pi_k = 1$:

$$F = \ln p(\boldsymbol{\mu}, \boldsymbol{\pi} | \{x_n\}_{n=1}^N) - \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

- So at first deriving for π_k affects two terms of the full posterior:

$$\frac{\partial F}{\partial \pi_k} = \frac{\sum_{n=1}^N \gamma(z_{nk})}{\pi_k} + \frac{\alpha_k - 1}{\pi_k} - \lambda = 0$$

Again, we can use $\sum_{n=1}^N \gamma(z_{nk}) = N_k$ to get:

$$\Rightarrow N_k + \alpha_k - 1 = \lambda \pi_k$$

The Lagrange Multiplier therefore is:

$$\begin{aligned}
\sum_{k=1}^K (N_k + \alpha_k - 1) &= \lambda \\
N + \sum_{k=1}^K \alpha_k - K &= \lambda
\end{aligned}$$

So, finally the update rule for π_k is:

$$\pi_k = \frac{N_k + \alpha_k - 1}{N + \sum_{k=1}^K \alpha_k - K}$$

- The second update rule is derived for μ_k :

$$\frac{\partial \ln p(\boldsymbol{\mu}, \boldsymbol{\pi} | \{x_n\}_{n=1}^N)}{\partial \mu_k} = \sum_{n=1}^N \gamma(z_{nk}) \left(\sum_{i=1}^D \frac{x_{ni}}{\mu_{ki}} - \frac{1 - x_{ni}}{1 - \mu_{ki}} \right) + \frac{a_k - 1}{\mu_k} - \frac{b_k - 1}{1 - \mu_k} = 0$$

Rearranging yields:

$$\sum_{n=1}^N \gamma(z_{nk}) \left(\frac{\mathbf{x}_n}{\mu_k} \right) + \frac{a_k - 1}{\mu_k} = \sum_{n=1}^N \gamma(z_{nk}) \left(\frac{1 - \mathbf{x}_n}{1 - \mu_k} \right) + \frac{b_k - 1}{1 - \mu_k}$$

Now, to get rid of the denominators, we multiply both μ_k and $(1 - \mu_k)$ on both sides:

$$(1 - \mu_k) \left(\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n + a_k - 1 \right) = \mu_k \left(\sum_{n=1}^N \gamma(z_{nk}) (1 - \mathbf{x}_n) + b_k - 1 \right)$$

Now, the term $\mu_{\mathbf{k}} \sum_{n=1}^N \gamma(z_{nk})$ cancels out

$$\rightarrow \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n + a_k - 1 = \mu_k \left(\sum_{n=1}^N \gamma(z_{nk}) + b_k - 1 + a_k - 1 \right)$$

Therefore, the M-Step update for $\mu_{\mathbf{k}}$ is:

$$\mu_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n + a_k - 1}{N_k + b_k + a_k - 2}$$