

# Machine Learning 2 - Homework 5

Andrii Skliar, 11636785

deadline: May 7, 2018

During the process of solving the homework problems, I have collaborated with the following colleagues:

Davide Belli   Gabriele Bani

*NB: credits for the Latex-format go to Iris Verweij, 2nd year MSc AI Student.*

**Problem 1:** Consider a Gaussian mixture model

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k N(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

1. Given the expected value of the complete-data log-likelihood (9.40 in Bishop's book)

$$\mathbb{E}_{\text{posterior}}[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{\ln \pi_k + \ln N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\}$$

Derive update rules for  $\boldsymbol{\pi}$ ,  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ .

**Solution:** To solve this problem, we need to first write down following Lagrangian:

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

Now, we can derive update rules for  $\boldsymbol{\pi}$ ,  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ .

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

$$\begin{aligned}
\frac{\partial F}{\partial \boldsymbol{\mu}_k} &= \{\text{using equation 86 from matrix cookbook}\} \\
&= \sum_{n=1}^N \gamma(z_{nk}) \left( -\frac{1}{2} \cdot -2 \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right) \\
&= \sum_{n=1}^N \gamma(z_{nk}) (\boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)) \\
&= \sum_{n=1}^N \gamma(z_{nk}) (\boldsymbol{\Sigma}_k^{-1} \mathbf{x}_n) - \sum_{n=1}^N \gamma(z_{nk}) (\boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k) = 0 \\
\Rightarrow N_k \cdot \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k &= \sum_{n=1}^N \gamma(z_{nk}) (\boldsymbol{\Sigma}_k^{-1} \mathbf{x}_n) \\
\boldsymbol{\mu}_k &= (N_k \cdot \boldsymbol{\Sigma}_k^{-1})^{-1} \sum_{n=1}^N \gamma(z_{nk}) (\boldsymbol{\Sigma}_k^{-1} \mathbf{x}_n) \\
\boldsymbol{\mu}_k &= \frac{1}{N_k} \boldsymbol{\Sigma}_k \boldsymbol{\Sigma}_k^{-1} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \\
\boldsymbol{\mu}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n
\end{aligned}$$

$$\begin{aligned}
\frac{\partial F}{\partial \Sigma_k} &= \{\text{using equations 57 and 61 from matrix cookbook} \\
&\quad \text{and the fact that } \Sigma \text{ is symmetric}\} \\
&= \sum_{n=1}^N -\frac{1}{2} \gamma(z_{nk}) \{ \Sigma_k^{-1} - \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} \} = 0 \\
-\frac{1}{2} \sum_{n=1}^N \gamma(z_{nk}) \Sigma_k^{-1} &= -\frac{1}{2} \sum_{n=1}^N \gamma(z_{nk}) \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} \\
N_k \cancel{\Sigma_k^{-1}} &= \Sigma_k^{-1} \cdot \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \cdot \cancel{\Sigma_k^{-1}} \\
\Sigma_k^{-1} &= N_k \left( \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \right)^{-1} \\
\Sigma_k &= \frac{\sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{N_k}
\end{aligned}$$

2. Consider a special case of the model above, in which the covariance matrices  $\Sigma_k$  of the components are all constrained to have a common value  $\Sigma$ . Derive EM equations for maximizing the likelihood function under such a model.

**Solution:**

$$\begin{aligned}
\frac{\partial F}{\partial \Sigma} &= -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{ \Sigma^{-1} - \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \Sigma^{-1} \} = 0 \\
\left( \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \right) \cancel{\Sigma^{-1}} &= \Sigma^{-1} \left( \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \right) \cancel{\Sigma^{-1}} \\
N &= \Sigma^{-1} \left( \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \right) \\
\Sigma^{-1} &= N \left( \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \right)^{-1} \\
\Sigma &= \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T
\end{aligned}$$

**Problem 2:** Suppose we wish to use the EM algorithm to maximize the posterior distribution  $p(\boldsymbol{\theta}|\mathbf{X})$  for a model containing latent variables  $\mathbf{z}$  and observed variables  $\mathbf{x}$ . Show that the E step remains the same as in the maximum likelihood case, where as in the M step, the quantity to be maximized is

$$\sum_{\mathbf{z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta})$$

**Solution:**

$$\begin{aligned} \ln p(\boldsymbol{\theta}|\mathbf{X}) &= \ln p(\boldsymbol{\theta}, \mathbf{X}) - \ln p(\mathbf{X}) \\ &= \ln p(\mathbf{X}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) - \ln p(\mathbf{X}) \\ &= \mathcal{L}(q, \boldsymbol{\theta}) + KL(q||p) + \ln p(\boldsymbol{\theta}) - \ln p(\mathbf{X}) \\ &\geq \mathcal{L}(q, \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) - \ln p(\mathbf{X}) \end{aligned}$$

Due to the fact, that  $q$  only appears in  $\mathcal{L}(q, \boldsymbol{\theta})$ , E-step will stay the same.  
M step update rule will be following:

$$\begin{aligned} \mathcal{L}(q, \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) - \ln p(\mathbf{X}) &= \mathcal{L}(q, \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) + \text{const} \\ &= \sum_{\mathbf{z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) + \sum_{\mathbf{z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\theta}^{old}) + \ln p(\boldsymbol{\theta}) + \text{const} \\ &= \sum_{\mathbf{z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) + \text{const} \end{aligned}$$

**Problem 3:** Derive the EM algorithm for maximizing the posterior probability  $p(\boldsymbol{\mu}, \boldsymbol{\pi}|\mathbf{x}_{n=1}^N)$  of Mixtures of Bernoulli distribution. (The E step is given in Bishops Book, you only need to do the M step).

**Solution:** Using equations derived in task 2, we can write equations for M step as following:

$$\begin{aligned} F &= \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] + \ln p(\boldsymbol{\theta}) \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{ \ln \pi_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln (1 - \mu_{ki})] \} \\ &\quad + \sum_{k=1}^K [\ln \text{Beta}(\boldsymbol{\mu}_k|a_k, b_k) + \ln \text{Dir}(\boldsymbol{\pi}_k|\boldsymbol{\alpha}_k)] \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{ \ln \pi_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln (1 - \mu_{ki})] \} \\ &\quad + \sum_{k=1}^K \left[ \sum_{i=1}^D \left\{ (a_k - 1) \ln \mu_{ki} + (b_k - 1) \ln (1 - \mu_{ki}) \right\} - \ln \mathcal{B}(a_k, b_k) \right. \\ &\quad \left. - \ln B(\boldsymbol{\alpha}_k) + (\alpha_k - 1) \ln \pi_k \right] \end{aligned}$$

Now we can perform M step for  $\mu$  and  $\pi$ :

$$\begin{aligned}
\frac{\partial F}{\partial \mu_{ki}} &= \sum_{n=1}^N \gamma(z_{nk}) \left( \frac{x_{ni}}{\mu_{ki}} - \frac{1-x_{ni}}{1-\mu_{ki}} \right) + \frac{a_k-1}{\mu_{ki}} - \frac{b_k-1}{1-\mu_{ki}} = 0 \\
\left[ \sum_{n=1}^N \gamma(z_{nk}) \frac{x_{ni}}{\mu_{ki}} \right] + \frac{a_k-1}{\mu_{ki}} &= \left[ \sum_{n=1}^N \gamma(z_{nk}) \frac{1-x_{ni}}{1-\mu_{ki}} \right] + \frac{b_k-1}{1-\mu_{ki}} \\
\frac{1}{\mu_{ki}} \left[ \sum_{n=1}^N \gamma(z_{nk}) x_{ni} + a_k - 1 \right] &= \frac{1}{1-\mu_{ki}} \left[ \sum_{n=1}^N \gamma(z_{nk}) (1-x_{ni}) + b_k - 1 \right] \\
\frac{1-\mu_{ki}}{\mu_{ki}} &= \frac{\sum_{n=1}^N \gamma(z_{nk}) (1-x_{ni}) + b_k - 1}{\sum_{n=1}^N \gamma(z_{nk}) x_{ni} + a_k - 1} \\
\frac{1}{\mu_{ki}} &= \frac{\sum_{n=1}^N \gamma(z_{nk}) (1-x_{ni}) + b_k - 1 + \sum_{n=1}^N \gamma(z_{nk}) x_{ni} + a_k - 1}{\sum_{n=1}^N \gamma(z_{nk}) x_{ni} + a_k - 1} \\
\frac{1}{\mu_{ki}} &= \frac{\sum_{n=1}^N \gamma(z_{nk}) + b_k + a_k - 2}{\sum_{n=1}^N \gamma(z_{nk}) x_{ni} + a_k - 1} \\
\mu_{ki} &= \frac{\sum_{n=1}^N \gamma(z_{nk}) x_{ni} + a_k - 1}{N_k + b_k + a_k - 2}
\end{aligned}$$

To perform M step for  $\pi$ , we also need to write proper Lagrangian, which in that case will be:

$$L = F + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

Therefore, M step update will look as following:

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