# Machine Learning 2 — Homework 5

Maurice Frank
11650656
maurice.frank@posteo.de

October 7, 2019

## Problem 1.

We have $p(\boldsymbol{x}) = \sum_{k=1}^{K} \pi_l \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$.

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

**1.**

For the update of $\pi$ we have to formulate the Lagrangian:

$$\mathcal{L} = \mathbb{E}_{\text{posterior}}\left[\ln p(\boldsymbol{X}, \boldsymbol{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})\right] + \lambda \cdot \left(\sum_{k=1}^{K} \pi_k - 1\right)$$

$$\frac{\partial}{\partial \pi_k}\mathcal{L} = \sum_{n=1}^{N} \frac{\partial}{\partial \pi_k}\gamma(z_{nk})\ln \pi_k + \lambda \cdot \frac{\partial}{\partial \pi_k}(\pi_k - 1)$$

$$= \sum_{n=1}^{N} \frac{\gamma(z_{nk})}{\pi_k} + \lambda$$

$$= \frac{1}{\pi_k} \cdot \left(\sum_{n=1}^{N} \gamma(z_{nk})\right) + \lambda$$

$$\equiv 0$$

$$\Longleftrightarrow$$

$$\pi_k = -\sum_{n=1}^{N} \frac{\gamma(z_{nk})}{\lambda}$$

$$= -\frac{N_k}{\lambda}$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \sum_{k=1}^{K} \pi_k - 1$$

$$\equiv 0 \quad \Longrightarrow$$

$$\sum_{k=1}^{K} \pi_k = 1 \quad \Longrightarrow$$

$$-\sum_{k=1}^{K} \frac{N_k}{\lambda} = 1 \quad \Longrightarrow$$

$$\lambda = -\sum_{k=1}^{K} N_k = N$$

$$\Longrightarrow$$

$$\pi_k = \frac{N_k}{N}$$

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

$$= \sum_{n=1}^{N} \gamma(z_{nk}) \frac{\frac{1}{2}(\boldsymbol{x}_n^T \boldsymbol{\Sigma}_k^{-1} + \boldsymbol{\Sigma}_k^{-1} \boldsymbol{x}_n - 2 \cdot \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k) \cdot \mathcal{N}(\boldsymbol{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))}{\mathcal{N}(\boldsymbol{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))}$$

$$= \sum_{n=1}^{N} \gamma(z_{nk}) \boldsymbol{\Sigma}_k^{-1} \cdot (\boldsymbol{x}_n - \boldsymbol{\mu}_k)$$

$$\equiv 0$$

$$\implies$$

$$\sum_{n=1}^{N} \gamma(z_{nk}) \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k = \sum_{n=1}^{N} \gamma(z_{nk}) \boldsymbol{\Sigma}_k^{-1} \boldsymbol{x}_n$$

$$\iff$$

$$\boldsymbol{\mu}_k = \frac{1}{\sum_{n=1}^{N} \gamma(z_{nk})} \sum_{n=1}^{N} \gamma(z_{nk}) \boldsymbol{x}_n$$

$$= \frac{\sum_{n=1}^{N} \gamma(z_{nk}) \boldsymbol{x}_n}{\sum_{n=1}^{N} \gamma(z_{nk})}$$

$$= \frac{\sum_{n=1}^{N} \gamma(z_{nk}) \boldsymbol{x}_n}{N_k}$$

$$\frac{\partial}{\partial \Sigma_k} \mathbb{E}_{\text{posterior}}[\ln p(X, Z | \mu, \Sigma, \pi)] = \sum_{n=1}^{N} \gamma(z_{nk}) \cdot \frac{\partial}{\partial \Sigma_k} \log \mathcal{N}(x_n | \mu_k, \Sigma_k)$$

$$= \sum_{n=1}^{N} \gamma(z_{nk}) \cdot \left( -\frac{1}{2} \frac{\partial}{\partial \Sigma_k} \ln |\Sigma_k| - \frac{1}{2} \frac{\partial}{\partial \Sigma_k} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right)$$

$$= \sum_{n=1}^{N} \gamma(z_{nk}) \cdot -\frac{1}{2} \left( \frac{|\Sigma_k| \Sigma_k^{-1}}{|\Sigma_k|} - \Sigma_k^{-1} (x_n - \mu_k)(x_n - \mu_k)^T \Sigma_k^{-1} \right)$$

$$= \sum_{n=1}^{N} \gamma(z_{nk}) \cdot \frac{1}{2} \left( \Sigma_k^{-1} (x_n - \mu_k)(x_n - \mu_k)^T \Sigma_k^{-1} - \Sigma_k^{-1} \right)$$

$$\equiv 0$$

$$\implies$$

$$\frac{N_k}{2} \Sigma_k^{-1} = \frac{1}{2} \sum_{n=1}^{N} \gamma(z_{nk}) \cdot \Sigma_k^{-1} (x_n - \mu_k)(x_n - \mu_k)^T \Sigma_k^{-1}$$

$$\implies$$

$$\Sigma_k = \sum_{n=1}^{N} \frac{\gamma(z_{nk})}{N_k} \cdot (x_n - \mu_k)(x_n - \mu_k)^T$$

**2.**

The updates given for $\pi_k$ and $\mu_k$ do not depend on the covariance. Therefore these two update rules do not change.

For the update of the common $\Sigma$ we have then:

$$\frac{\partial}{\partial \Sigma} \mathbb{E}_{\text{posterior}}[\ln p(X, Z | \mu, \Sigma, \pi)] = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \cdot \frac{\partial}{\partial \Sigma} \log \mathcal{N}(x_n | \mu_k, \Sigma)$$

$$= \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \cdot \frac{1}{2} \left( \Sigma^{-1} (x_n - \mu_k)(x_n - \mu_k)^T \Sigma^{-1} - \Sigma^{-1} \right)$$

$$\equiv 0$$

$$\implies$$

$$\Sigma = \sum_{n=1}^{N} \sum_{k=1}^{K} y \frac{\gamma(z_{nk})}{N_k} \cdot (x_n - \mu_k)(x_n - \mu_k)^T$$

# Problem 2.

$$\ln p(\boldsymbol{\theta}|\boldsymbol{X}) = \ln p(\boldsymbol{\theta}, \boldsymbol{X}) - \ln p(\boldsymbol{X})$$
$$= \ln p(\boldsymbol{X}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) - \ln p(\boldsymbol{X})$$
$$= \mathcal{L}(q, \boldsymbol{\theta}) + \mathrm{KL}[q||p] + \ln p(\boldsymbol{\theta}) - \ln p(\boldsymbol{X})$$
$$\geqslant \mathcal{L}(q, \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) - \ln p(\boldsymbol{X})$$

For the E-Step we maximize the lower bound w.r.t. $q$. As only $\mathcal{L}(q, \boldsymbol{\theta})$ is dependent on $q$ we have the same situation as in the E-step for the ML estimate.

For the M-step we maximize the lower bound w.r.t. $\boldsymbol{\theta}$:

$$\arg\max_{\boldsymbol{\theta}} \mathcal{L}(q, \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) - \ln p(\boldsymbol{X}) = \arg\max_{\boldsymbol{\theta}} \mathcal{L}(q, \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta})$$

$$= \arg\max_{\boldsymbol{\theta}} \sum_{z} p(\boldsymbol{Z}|\boldsymbol{X}, \boldsymbol{\theta}^{\mathrm{old}}) \cdot \ln \frac{p(\boldsymbol{X}, \boldsymbol{Z}|\boldsymbol{\theta})}{p(\boldsymbol{Z}|\boldsymbol{X}, \boldsymbol{\theta}^{\mathrm{old}})} + \ln p(\boldsymbol{\theta})$$

$$= \arg\max_{\boldsymbol{\theta}} \sum_{z} p(\boldsymbol{Z}|\boldsymbol{X}, \boldsymbol{\theta}^{\mathrm{old}}) \cdot \ln p(\boldsymbol{X}, \boldsymbol{Z}|\boldsymbol{\theta})$$
$$- \sum_{z} p(\boldsymbol{Z}|\boldsymbol{X}, \boldsymbol{\theta}^{\mathrm{old}}) \cdot \ln p(\boldsymbol{Z}|\boldsymbol{X}, \boldsymbol{\theta}^{\mathrm{old}}) + \ln p(\boldsymbol{\theta})$$

$$= \arg\max_{\boldsymbol{\theta}} \sum_{z} p(\boldsymbol{Z}|\boldsymbol{X}, \boldsymbol{\theta}^{\mathrm{old}}) \ln p(\boldsymbol{X}, \boldsymbol{Z}|\boldsymbol{\theta}) + H[p(\boldsymbol{Z}|\boldsymbol{X}, \boldsymbol{\theta}^{\mathrm{old}})] + \ln p(\boldsymbol{\theta})$$

$$= \arg\max_{\boldsymbol{\theta}} \sum_{z} p(\boldsymbol{Z}|\boldsymbol{X}, \boldsymbol{\theta}^{\mathrm{old}}) \ln p(\boldsymbol{X}, \boldsymbol{Z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta})$$

The entropy of $p(\boldsymbol{Z}|\boldsymbol{X}, \boldsymbol{\theta}^{\mathrm{old}})$ drops out as it does not depend on $\boldsymbol{\theta}$.

# Problem 3.

We have:

$$\boldsymbol{\pi}|\boldsymbol{\alpha} \sim \mathrm{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha})$$
$$\boldsymbol{z}_n|\boldsymbol{\pi} \sim \mathrm{Mult}(\boldsymbol{z}_n|\boldsymbol{\pi})$$
$$\mu_k|a_k, b_k \sim \mathrm{Beta}(\mu|a_k, b_k)$$
$$\boldsymbol{x}_n|\boldsymbol{z}_n, \boldsymbol{\mu} = \{\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_K\} \sim \prod_{k=1}^{K} (\mathrm{Bern}(\boldsymbol{x}_n|\boldsymbol{\mu}_k))^{z_{nk}}$$

Further lets write some stuff that we gonna use later:

$$\ln p(\boldsymbol{\mu}) = \sum_{k=1}^{K} \sum_{j}^{D} (a_k - 1) \ln \mu_{kj} + (b_k - 1) \ln(1 - \mu_{kj}) - \ln B(a_k, b_k)$$

$$\ln p(\boldsymbol{\pi}) = \sum_{k=1}^{K} (\alpha_k - 1) \ln \pi_k - \ln B(\alpha_k)$$

$$\mathbb{E}_{p(\boldsymbol{Z}|\boldsymbol{X},\boldsymbol{\mu}^{\text{old}},\boldsymbol{\pi}^{\text{old}})}[ln p(\boldsymbol{X}, \boldsymbol{Z}|\boldsymbol{\mu}, \boldsymbol{\pi})] = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \left\{ \ln \pi_k + \sum_{i=1}^{D} [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki})] \right\}$$

with the last being already the result from the E-step as described in the book.

To calculate the updates for $\boldsymbol{\mu}$ and $\boldsymbol{\pi}$ in the M-step we can use the result from Problem 2:

$$\underset{\boldsymbol{\mu},\boldsymbol{\pi}}{\arg\max} \, \mathcal{L}(q, \boldsymbol{\mu}, \boldsymbol{\pi}) + \ln p(\boldsymbol{\theta}) + \ln p(\boldsymbol{\pi}) = \underset{\boldsymbol{\mu},\boldsymbol{\pi}}{\arg\max} \sum_{z} p(\boldsymbol{Z}|\boldsymbol{X}, \boldsymbol{\mu}^{\text{old}}, \boldsymbol{\pi}^{\text{old}}) \ln p(\boldsymbol{X}, \boldsymbol{Z}|\boldsymbol{\mu}, \boldsymbol{\pi}) + \ln p(\boldsymbol{\mu}) + \ln p(\boldsymbol{\pi})$$

$$= \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \left\{ \ln \pi_k + \sum_{i=1}^{D} [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki})] \right\}$$

$$+ \sum_{k=1}^{K} \sum_{j}^{D} (a_k - 1) \ln \mu_{kj} + (b_k - 1) \ln(1 - \mu_{kj}) - \ln B(a_k, b_k)$$

$$+ \sum_{k=1}^{K} (\alpha_k - 1) \ln \pi_k - \ln B(\alpha_k)$$

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

we already added the lagrange Multiplier for the condition of the Dirichlet distribution of $\pi$.

Now we maximize for the two variables.

First for $\mu$:

$$\frac{\partial}{\partial \mu_{ki}} \mathcal{L}(q, \mu, \pi) + \ln p(\theta) + \ln p(\pi) = \sum_{n=1}^{N} \gamma(z_{nk}) \left\{ x_{ni} \frac{1}{\mu_{ki}} - (1 - x_{ni}) \frac{1}{1 - \mu_{ki}} \right\}$$

$$+ (a_k - 1) \frac{1}{\mu_{ki}} - (b_k - 1) \frac{1}{1 - \mu_{ki}}$$

$$= \frac{1}{\mu_{ki}} \left( \sum_{n=1}^{N} \gamma(z_{nk}) x_{ni} + (a_k - 1) \right)$$

$$- \frac{1}{1 - \mu_{ki}} \left( \sum_{n=1}^{N} \gamma(z_{nk})(1 - x_{ni}) + (b_k - 1) \right)$$

$$\equiv 0$$

$$\Longrightarrow$$

$$\left( \sum_{n=1}^{N} \gamma(z_{nk}) x_{ni} + (a_k - 1) \right) = \mu_{ki} \left[ \left( \sum_{n=1}^{N} \gamma(z_{nk})(1 - x_{ni}) + (b_k - 1) \right) + \left( \sum_{n=1}^{N} \gamma(z_{nk}) x_{ni} + (a_k - 1) \right) \right]$$

$$\mu_{ki} = \frac{\sum_{n=1}^{N} \gamma(z_{nk}) x_{ni} + (a_k - 1)}{\sum_{n=1}^{N} \gamma(z_{nk}) + b_k + a_k - 2}$$

$$= \frac{\sum_{n=1}^{N} \gamma(z_{nk}) x_{ni} + (a_k - 1)}{N_k + b_k + a_k - 2}$$

And lastly for $\boldsymbol{\pi}$:

$$\frac{\partial}{\partial \pi_k} \mathcal{L}(q, \boldsymbol{\mu}, \boldsymbol{\pi}) + \ln p(\boldsymbol{\theta}) + \ln p(\boldsymbol{\pi}) = \sum_{n=1}^{N} \gamma(z_{nk}) \frac{1}{\pi_k}$$

$$+ (\alpha_k - 1) \frac{1}{\pi_k}$$

$$+ \lambda$$

$$= \frac{1}{\pi_k} \cdot (N_k + \alpha_k - 1) + \lambda$$

$$\equiv 0$$

$$\implies$$

$$\pi_k = -\frac{N_k + \alpha_k - 1}{\lambda}$$

and:

$$\lambda = \lambda \sum_{k=1}^{K} \pi_k$$

$$= -\left( \sum_{k=1}^{K} N_k + \sum_{k=1}^{K} \alpha_k - \sum_{k=1}^{K} 1 \right)$$

$$= -\left( N + \sum_{k=1}^{K} \alpha_k - K \right)$$

$$\implies$$

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

8