

## Machine Learning Exercise Sheet 12

### Clustering

### In-class Exercises

#### K-Medians

**Problem 1:** Consider a modified version of the  $K$ -means objective, where we use  $L_1$  distance instead.

$$\mathcal{J}(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}) = \sum_{i=1}^N \sum_{k=1}^K z_{ik} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_1$$

This variation of the algorithm is called  $K$ -medians. Derive the Lloyd's algorithm for this model.

1. Updating the cluster assignments  $z_{ik}$  is the same as for the  $K$ -means algorithm:

$$z_{ik}^{new} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_i - \boldsymbol{\mu}_j\|_1 \\ 0 & \text{else.} \end{cases}$$

2. The updates for  $\boldsymbol{\mu}_k$ 's should solve

$$\boldsymbol{\mu}_k^{new} = \arg \min_{\boldsymbol{\mu}_k} \sum_{i=1}^N z_{ik} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_1$$

The objective for each single centroid  $\boldsymbol{\mu}_k$  can be rewritten as

$$\begin{aligned} \mathcal{J}(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}_k) &= \sum_{i=1}^N z_{ik} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_1 \\ &= \sum_{i=1}^N z_{ik} \sum_{d=1}^D |\mathbf{x}_{id} - \boldsymbol{\mu}_{kd}| \end{aligned}$$

Clearly, this is a convex function of  $\boldsymbol{\mu}_k$ , as it is a sum of piecewise linear functions. We can actually solve for each  $\boldsymbol{\mu}_{kd}$  separately, as they do not interact in the objective, by finding the roots of the derivatives.

Observe, that

$$\frac{\partial}{\partial \boldsymbol{\mu}_{kd}} |\mathbf{x}_{id} - \boldsymbol{\mu}_{kd}| = \begin{cases} \frac{\partial}{\partial \boldsymbol{\mu}_{kd}} (\boldsymbol{\mu}_{kd} - \mathbf{x}_{id}) = 1 & \text{if } \boldsymbol{\mu}_{kd} > \mathbf{x}_{id} \\ \frac{\partial}{\partial \boldsymbol{\mu}_{kd}} (\mathbf{x}_{id} - \boldsymbol{\mu}_{kd}) = -1 & \text{if } \boldsymbol{\mu}_{kd} < \mathbf{x}_{id} \\ 0 & \text{if } \boldsymbol{\mu}_{kd} = \mathbf{x}_{id}. \end{cases}$$

(Note: actually the absolute value function is not differentiable at 0, so the derivative is undefined. A rigorous treatment of this problem would require us to use subgradients (see [https://web.stanford.edu/class/ee364b/lectures/subgradients\\_notes.pdf](https://web.stanford.edu/class/ee364b/lectures/subgradients_notes.pdf)), but just "pretending" that the gradient is 0 suffices for our purpose.)

Hence, the derivative of the entire objective is

$$\begin{aligned}\frac{\partial}{\partial \mu_{kd}} \mathcal{J}(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}) &= \sum_{i=1}^N z_{ik} |x_{id} - \mu_{kd}| \\ &= \sum_{i=1}^N z_{ik} \mathbb{I}[\mu_{kd} > x_{id}] - \sum_{i=1}^N z_{ik} \mathbb{I}[\mu_{kd} < x_{id}] \stackrel{!}{=} 0\end{aligned}$$

The first sum represents "number of points  $\mathbf{x}_i$  assigned to class  $k$ , such that  $x_{id} < \mu_{kd}$ ". Each of these sums represents the number of points in class  $k$ , that are located to the left (right) of the given value of  $\mu_{kd}$ . Because we want to set the gradient to zero, we are looking for such a  $\mu_{kd}$ , that along the axis  $d$  exactly  $N_k/2$  points are to left of it, and another  $N_k/2$  points are to the right (where  $N_k = \sum_{i=1}^N z_{ik}$ ). This is exactly the definition of a *median*.

Therefore, the optimal update is given as

$$\mu_{kd} = \text{median} \{x_{id} \text{ such that } z_{ik} = 1\}$$

## Gaussian Mixture Model

**Problem 2:** Derive the E-step update for the Gaussian mixture model.

In the E-step we have to evaluate the posterior distribution over the latent variables given the current parameters, i.e.  $\gamma_t(\mathbf{Z})$ . Because GMMs assume that the latent variables are independent,  $\gamma_t(\mathbf{Z}) = \prod_{i=1}^N \gamma_t(\mathbf{z}_i)$  and it is enough to derive the E-step for a single data point. The update rule follows directly from Bayes' theorem.

$$\begin{aligned}\gamma_t(\mathbf{z}_i = k) &= p(\mathbf{z}_i = k \mid \mathbf{x}_i, \boldsymbol{\pi}^{(t)}, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}) \\ &= \frac{p(\mathbf{x}_i \mid \mathbf{z}_i = k, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}) p(\mathbf{z}_i = k \mid \boldsymbol{\pi}^{(t)})}{p(\mathbf{x}_i \mid \boldsymbol{\pi}^{(t)}, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)})} \\ &= \frac{p(\mathbf{x}_i \mid \mathbf{z}_i = k, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}) p(\mathbf{z}_i = k \mid \boldsymbol{\pi}^{(t)})}{\sum_{j=1}^K p(\mathbf{x}_i \mid \mathbf{z}_i = j, \boldsymbol{\pi}^{(t)}, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}) p(\mathbf{z}_i = j \mid \boldsymbol{\pi}^{(t)})} \\ &= \frac{\pi_k^{(t)} \mathcal{N}(\mathbf{x}_i \mid \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)})}{\sum_{j=1}^K \pi_j^{(t)} \mathcal{N}(\mathbf{x}_i \mid \boldsymbol{\mu}_j^{(t)}, \boldsymbol{\Sigma}_j^{(t)})}\end{aligned}$$

**Problem 3:** Derive the M-step update for the Gaussian mixture model.

In the M-step we maximize  $\mathcal{L} = \mathbb{E}_{\mathbf{Z} \sim \gamma_t(\mathbf{Z})} [\log p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})]$  with respect to  $\boldsymbol{\pi}$ ,  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . When we plug in the definition of the expected value and expand, we get

$$\begin{aligned} \mathcal{L} &= \sum_{i=1}^N \sum_{k=1}^K \gamma_t(\mathbf{z}_i = k) \log p(\mathbf{x}_i, \mathbf{z}_i = k \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= \sum_{i=1}^N \sum_{k=1}^K \gamma_t(\mathbf{z}_i = k) \log p(\mathbf{x}_i \mid \mathbf{z}_i = k, \boldsymbol{\mu}, \boldsymbol{\Sigma}) p(\mathbf{z}_i = k \mid \boldsymbol{\pi}) \\ &= \underbrace{\sum_{i=1}^N \sum_{k=1}^K \gamma_t(\mathbf{z}_i = k) \log p(\mathbf{x}_i \mid \mathbf{z}_i = k, \boldsymbol{\mu}, \boldsymbol{\Sigma})}_{\mathcal{L}_x} + \underbrace{\sum_{i=1}^N \sum_{k=1}^K \gamma_t(\mathbf{z}_i = k) \log p(\mathbf{z}_i = k \mid \boldsymbol{\pi})}_{\mathcal{L}_z} \end{aligned}$$

where  $\mathcal{L}_z$  only depends on  $\boldsymbol{\pi}$  and  $\mathcal{L}_x$  only depends on  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . To find the optimal  $\boldsymbol{\pi}$ , we need to maximize  $\mathcal{L}_z$  with respect to  $\boldsymbol{\pi}$ . Since  $\boldsymbol{\pi}$  has several constraints placed on it, we will have to solve the following convex optimization problem.

$$\begin{aligned} &\text{maximize} \quad \mathcal{L}_z \\ &\text{subject to} \quad \sum_{k=1}^K \pi_k - 1 = 0 \end{aligned}$$

Before we formulate the Lagrangian, we simplify  $\mathcal{L}_z$  as

$$\mathcal{L}_z = \sum_{i=1}^N \sum_{k=1}^K \gamma_t(\mathbf{z}_i = k) \log p(\mathbf{z}_i = k \mid \boldsymbol{\pi}) = \sum_{k=1}^K N_k \log \pi_k$$

where  $N_k = \sum_{i=1}^N \gamma_t(\mathbf{z}_i = k)$  is the size of the  $k$ -th cluster. The Lagrangian is given by

$$f(\boldsymbol{\pi}, \lambda) = \sum_{k=1}^K N_k \log \pi_k + \lambda \left( 1 - \sum_{k=1}^K \pi_k \right)$$

and it has its maximum in  $\boldsymbol{\pi}$  at

$$\frac{\partial f}{\partial \pi_k} = \frac{N_k}{\pi_k} - \lambda \stackrel{!}{=} 0 \Leftrightarrow \pi_k = \frac{N_k}{\lambda}$$

because  $f$  is concave as a function of  $\boldsymbol{\pi}$ . This gives us the dual function as

$$g(\lambda) = \max_{\boldsymbol{\pi}} f(\boldsymbol{\pi}, \lambda) = f\left(\left(\frac{N_1}{\lambda}, \dots, \frac{N_K}{\lambda}\right), \lambda\right) = \sum_{k=1}^K N_k \log \frac{N_k}{\lambda} + \lambda - N.$$

When  $f$  is concave, the dual is convex and we find the minimum of  $g$  at

$$\frac{\partial g}{\partial \lambda} = \sum_{k=1}^K N_k \frac{\lambda}{N_k} \left( -\frac{N_k}{\lambda^2} \right) + 1 = 1 - \frac{N}{\lambda} \stackrel{!}{=} 0 \Leftrightarrow \lambda = N.$$

This means that the M-step for  $\boldsymbol{\pi}$  is  $\boldsymbol{\pi}_k^{(t+1)} = \frac{N_k}{N}$ .

To find the M-step rules for  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ , we need to examine  $\mathcal{L}_{\mathbf{x}}$ .

$$\begin{aligned}\mathcal{L}_{\mathbf{x}} &= \sum_{i=1}^N \sum_{k=1}^K \gamma_t(z_i = k) \log(\mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \\ &= -\frac{1}{2} \sum_{i=1}^N \sum_{k=1}^K \gamma_t(z_i = k) \left( (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) + D \log(2\pi) + \log \det \boldsymbol{\Sigma}_k \right).\end{aligned}$$

where  $D$  is the feature dimension. We can take the derivative with respect to  $\boldsymbol{\mu}_k$

$$\frac{\partial \mathcal{L}_{\mathbf{x}}}{\partial \boldsymbol{\mu}_k} = -\frac{1}{2} \sum_{i=1}^N \gamma_t(z_i = k) \left( (-1) \cdot \left( \boldsymbol{\Sigma}_k^{-1} + \boldsymbol{\Sigma}_k^{-T} \right) (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) = \sum_{i=1}^N \gamma_t(z_i = k) \left( \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right)$$

and then find its root

$$\frac{\partial \mathcal{L}_{\mathbf{x}}}{\partial \boldsymbol{\mu}_k} = 0 \Leftrightarrow \sum_{i=1}^N \gamma_t(z_i = k) \boldsymbol{\Sigma}_k^{-1} \mathbf{x}_i = N_k \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k \Leftrightarrow \boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i=1}^N \gamma_t(z_i = k) \mathbf{x}_i$$

which gives us the update rule

$$\boldsymbol{\mu}_k^{(t+1)} = \frac{1}{N_k} \sum_{i=1}^N \gamma_t(z_i = k) \mathbf{x}_i.$$

It remains to find the M-step for  $\boldsymbol{\Sigma}$ . Again we proceed by taking the derivative with respect to  $\boldsymbol{\Sigma}_k$

$$\begin{aligned}\frac{\partial \mathcal{L}_{\mathbf{x}}}{\partial \boldsymbol{\Sigma}_k} &= -\frac{1}{2} \sum_{i=1}^N \gamma_t(z_i = k) \left[ -\boldsymbol{\Sigma}_k^{-T} (\mathbf{x}_i - \boldsymbol{\mu}_k) (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-T} + \boldsymbol{\Sigma}_k^{-T} \right] \\ &= -\frac{1}{2} \left( N_k I_D - \sum_{i=1}^N \gamma_t(z_i = k) \left[ \boldsymbol{\Sigma}_k^{-T} (\mathbf{x}_i - \boldsymbol{\mu}_k) (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \right] \right) \boldsymbol{\Sigma}_k^{-T}\end{aligned}$$

where  $I_D$  is the  $D$ -dimensional identity matrix. We finish by finding its root

$$\frac{\partial \mathcal{L}_{\mathbf{x}}}{\partial \boldsymbol{\Sigma}_k} = 0 \Leftrightarrow N_k I_D = \boldsymbol{\Sigma}_k^{-T} \sum_{i=1}^N \gamma_t(z_i = k) (\mathbf{x}_i - \boldsymbol{\mu}_k) (\mathbf{x}_i - \boldsymbol{\mu}_k)^T$$

which produces the final update rule

$$\boldsymbol{\Sigma}_k^{(t+1)} = \frac{1}{N_k} \sum_{i=1}^N \gamma_t(z_i = k) (\mathbf{x}_i - \boldsymbol{\mu}_k) (\mathbf{x}_i - \boldsymbol{\mu}_k)^T.$$

In this exercise we have used the following matrix calculus rules which you can look up in the matrix cookbook.

$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{a}} = (\mathbf{X} + \mathbf{X}^T) \mathbf{a}^T \quad \frac{\partial \mathbf{a}^T \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -\mathbf{X}^{-T} \mathbf{b} \mathbf{a}^T \mathbf{X}^{-T} \quad \frac{\partial \log |\det \mathbf{X}|}{\partial \mathbf{X}} = \mathbf{X}^{-T}$$

## Expectation Maximization Algorithm

**Problem 4:** Consider a mixture model where the components are given by independent Bernoulli variables. This is useful when modelling, e.g., binary images, where each of the  $D$  dimensions of the image  $\mathbf{x}$  corresponds to a different pixel that is either black or white. More formally, we have

$$p(\mathbf{x} \mid \mathbf{z} = k) = \prod_{d=1}^D \theta_{kd}^{x_d} (1 - \theta_{kd})^{1-x_d}.$$

That is, for a given mixture index  $\mathbf{z} = k$ , we have a product of independent Bernoullis, where  $\theta_{kd}$  denotes the Bernoulli parameter for component  $k$  at pixel  $d$ .

Derive the EM algorithm for the parameters  $\boldsymbol{\theta} = \{\theta_{kd} \mid k = 1, \dots, K, d = 1, \dots, D\}$  of a mixture of Bernoullis.

Assume here for simplicity, that the distribution of components  $p(\mathbf{z})$  is uniform:  $p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} = \prod_{k=1}^K \left(\frac{1}{K}\right)^{z_k}$ .

Due to the uniform prior on  $\mathbf{z}_i$ , the  $p(\mathbf{z}_i)$  cancel and the responsibilities compute as

$$\gamma_t(\mathbf{z}_i = k) = \frac{p(\mathbf{x}_i \mid \mathbf{z}_i = k, \boldsymbol{\theta}) \cdot p(\mathbf{z}_i = k)}{\sum_{l=1}^K p(\mathbf{x}_i \mid \mathbf{z}_i = l, \boldsymbol{\theta}) \cdot p(\mathbf{z}_i = l)} = \frac{p(\mathbf{x}_i \mid \mathbf{z}_i = k, \boldsymbol{\theta})}{\sum_{l=1}^K p(\mathbf{x}_i \mid \mathbf{z}_i = l, \boldsymbol{\theta})}$$

which constitutes the E-step.

It remains to derive the M-step. Similar to mixture of Gaussians:

$$\begin{aligned} \mathbb{E}_{\mathbf{z} \sim \gamma_t(\mathbf{z})} [\log p(\mathbf{X}, \mathbf{z} \mid \boldsymbol{\theta}^{(t)})] &= \sum_{i=1}^N \sum_{k=1}^K \gamma_t(\mathbf{z}_i = k) \log \left( \frac{1}{K} \prod_{d=1}^D \theta_{kd}^{x_{id}} (1 - \theta_{kd})^{1-x_{id}} \right) \\ &= C + \underbrace{\sum_{i=1}^N \sum_{k=1}^K \gamma_t(\mathbf{z}_i = k) \sum_{d=1}^D (x_{id} \log \theta_{kd} + (1 - x_{id}) \log(1 - \theta_{kd}))}_{=: \mathcal{L}_i} \end{aligned}$$

The constant  $C$  collects all terms independent of  $\boldsymbol{\theta}$  and hence irrelevant for further optimization.

We now need to take derivatives with respect to  $\boldsymbol{\theta}$ .

$$\begin{aligned} \frac{\partial \mathcal{L}_i}{\partial \theta_{k',d'}} &= \sum_{k=1}^K \gamma_t(\mathbf{z}_i = k) \sum_{d=1}^D \left( x_{id} \frac{\partial \log \theta_{kd}}{\partial \theta_{k',d'}} + (1 - x_{id}) \frac{\partial \log(1 - \theta_{kd})}{\partial \theta_{k',d'}} \right) \\ &= \gamma_t(\mathbf{z}_i = k) \left( \frac{x_{id}}{\theta_{k',d'}} - \frac{1 - x_{id}}{1 - \theta_{k',d'}} \right) \end{aligned}$$

We observe that the  $\theta_{kd}$  do not interact, so their optimal values are independent from each other and we can handle them individually.

$$\frac{\partial \mathbb{E}_{\mathbf{z} \sim \gamma_t(\mathbf{z})} [\log p(\mathbf{X}, \mathbf{z} \mid \boldsymbol{\theta})]}{\partial \theta_{kd}} = \sum_{i=1}^N \frac{\partial \mathcal{L}_i}{\partial \theta_{kd}} = \sum_{i=1}^N \gamma_t(\mathbf{z}_i = k) \left( \frac{x_{id}}{\theta_{kd}} - \frac{1 - x_{id}}{1 - \theta_{kd}} \right)$$

By finding the roots  $\frac{\partial \mathbb{E}_{\mathbf{z} \sim \gamma_t(\mathbf{z})} [\log p(\mathbf{X}, \mathbf{z} | \boldsymbol{\theta})]}{\partial \boldsymbol{\theta}_{kd}} = 0$ , we obtain the optimal update in a similar fashion as in the standard Bernoulli MLE:

$$\boldsymbol{\theta}_{kd} = \frac{\sum_{i=1}^N \gamma_t(\mathbf{z}_i = k) \mathbf{x}_{id}}{\sum_{i=1}^N \gamma_t(\mathbf{z}_i = k)}$$