

Machine Learning Exercise Sheet 12

Clustering

Exercise sheets consist of two parts: In-class exercises and homework. The in-class exercises will be solved and discussed during the tutorial. The homework is for you to solve at home and further engage with the lecture content. There is no grade bonus and you do not have to upload any solutions. Note that the order of some exercises might have changed compared to last year's recordings.

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.

Gaussian Mixture Model

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

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

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 θ_{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: $p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} = \prod_{k=1}^K \left(\frac{1}{K}\right)^{z_k}$.

P₁
$$J(X, Z, \mu) = \sum_{i=1}^N \sum_{k=1}^K z_{ik} \|x_i - \mu_k\|$$

$$\min_z J(X, Z, \mu) \quad z_{ik} = \begin{cases} 1, & \text{if } k = \arg \min_j \|x_i - \mu_j\| \\ 0, & \text{else} \end{cases}$$

$$\min_{\mu} J(X, Z, \mu) \quad J(X, Z, \mu) = \sum_{i=1}^N z_{ik} \sum_{k=1}^K \|x_{id} - \mu_{kd}\|$$

convex function

$$\frac{\partial \|x_{id} - \mu_{kd}\|}{\partial \mu_{kd}} = \begin{cases} \frac{\partial \mu_{kd} - x_{id}}{\partial \mu_{kd}} = 1 & \mu_{kd} > x_{id} \\ \frac{\partial x_{id} - \mu_{kd}}{\partial \mu_{kd}} = -1 & \mu_{kd} < x_{id} \\ 0 & \mu_{kd} = x_{id} \end{cases}$$

$$\frac{\partial J(X, Z, \mu)}{\partial \mu_{kd}} = \sum_{i=1}^N z_{ik} \cdot \mathbb{I}(\mu_{kd} > x_{id}) - \sum_{i=1}^N z_{ik} \mathbb{I}(\mu_{kd} < x_{id}) \stackrel{!}{=} 0$$

$\frac{N_k}{2}$ points $\quad \quad \quad \frac{N_k}{2}$ points

$$N_k = \sum_{i=1}^N z_{ik}$$

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

P₂
$$\gamma_0(z) = p(z | X, \tau^{(0)}, \mu^{(0)}, \Sigma^{(0)})$$

$$\begin{aligned} \gamma_t(z = z_{ik}) &= p(z_{ik} | x_i, \tau^{(t)}, \mu^{(t)}, \Sigma^{(t)}) = \frac{p(x_i | z_{ik}, \mu^{(t)}, \Sigma^{(t)}) \cdot p(z_{ik} | \tau^{(t)})}{p(x_i | \tau^{(t)}, \mu^{(t)}, \Sigma^{(t)})} \\ &= \frac{\tau_k^{(t)} \cdot \mathcal{N}(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum p(x_i, z_{ik} | \tau^{(t)}, \mu^{(t)}, \Sigma^{(t)})} \\ &= \frac{\tau_k^{(t)} \cdot \mathcal{N}(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum p(x_i | z = z_{ik}, \tau^{(t)}, \mu^{(t)}, \Sigma^{(t)}) \cdot p(z = z_{ik} | \tau^{(t)})} \\ &= \frac{\tau_k^{(t)} \mathcal{N}(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_j^K \tau_j^{(t)} \mathcal{N}(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})} \end{aligned}$$

P₃
$$\mu^{t+1}, \Sigma^{t+1}, \tau^{t+1} = \underset{\mu, \Sigma, \tau}{\operatorname{argmax}} \bar{E}[\log p(X, Z | \tau, \mu, \Sigma)]$$

$$= \underset{\mu, \Sigma, \tau}{\operatorname{argmax}} \sum \gamma_t(z) \cdot \log p(x, z | \tau, \mu, \Sigma)$$

$$= \sum_{i=1}^N \sum_{k=1}^K \gamma_t(z_i) \cdot \log p(x_i | z_i = k, \tau, \mu, \Sigma) \cdot p(z_i = k | \tau)$$

$$\begin{aligned} &= \underbrace{\sum_{i=1}^N \sum_{k=1}^K \gamma_t(z_i) / n \cdot p(x_i | z_i = k, \tau, \mu, \Sigma)}_{\sim \mu, \Sigma} + \underbrace{\sum_{i=1}^N \sum_{k=1}^K \gamma_t(z_i) / n \cdot p(z_i = k | \tau)}_{\sim \tau} \\ &= \sum_{i=1}^N \sum_{k=1}^K \gamma_t(z_i) \cdot \log \tau_k \end{aligned}$$

$$\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) N_k \bar{z}_k$$

$$\sum_{i=1}^n \sum_{k=1}^K \gamma_t(z_i=k) \ln p(x_i | z_i=k, \mu, \Sigma)$$

$$= \sum_{i=1}^n \sum_{k=1}^K \gamma_t(z_i=k) \cdot \left[\ln \frac{1}{\sqrt{2\pi}} + \ln \frac{1}{\sqrt{|\Sigma_k|}} - \frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right]$$

$$\left[-\frac{1}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right]$$

$$\frac{\partial \mathcal{L}}{\partial \mu_k} = -\frac{1}{2} \sum_{i=1}^n \gamma_t(z_i=k) (-2 x_i^T \Sigma_k^{-1} + 2 \Sigma_k^{-1} \mu_k)$$

$$= \sum_{i=1}^n \gamma_t(z_i=k) \Sigma_k^{-1} (x_i - \mu_k) \stackrel{!}{=} 0$$

$$\sum_{i=1}^n \gamma_t(z_i=k) \Sigma_k^{-1} x_i^T = N_k \Sigma_k^{-1} \mu_k$$

$$\mu_k = \frac{1}{N_k} \sum_{i=1}^n \gamma_t(z_i=k) x_i^T$$

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

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \frac{N_k}{\pi_k} - \lambda \stackrel{!}{=} 0$$

$$\pi_k = \frac{1}{\lambda} N_k$$

$$g(x) = \sum_{k=1}^K N_k \cdot \ln \frac{N_k}{\lambda} + \lambda \left(1 - \sum_{k=1}^K \frac{N_k}{\lambda}\right)$$

$$= \sum_{k=1}^K N_k \cdot \ln \frac{N_k}{\lambda} + \lambda - N$$

$$\frac{\partial g}{\partial \lambda} = -\sum_{k=1}^K N_k \cdot \frac{1}{\lambda} + 1 \stackrel{!}{=} 0$$

$$= -\frac{N}{\lambda} + 1$$

$$\lambda = N$$

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

$$\frac{\partial \mathcal{L}}{\partial \Sigma_k} = -\frac{1}{2} \sum_{i=1}^n \gamma_t(z_i=k) \left[\Sigma_k^{-1} + (-\Sigma_k^{-1} x_i x_i^T \Sigma_k^{-1}) - 2(-\Sigma_k^{-1} \mu_k x_i^T \Sigma_k^{-1}) + (-\Sigma_k^{-1} \mu_k \mu_k^T \Sigma_k^{-1}) \right]$$

$$\Sigma_k^{-1} - \Sigma_k^{-1} \left[x_i x_i^T - 2 \mu_k x_i^T + \mu_k \mu_k^T \right]$$

$$\Sigma_k^{-1} - \Sigma_k^{-1} \left[(x_i - \mu_k) (x_i - \mu_k)^T \right] \Sigma_k^{-1}$$

$$= -\frac{1}{2} \left[N_k \cdot \Sigma_k^{-1} - \sum_{i=1}^n \gamma_t(z_i=k) \Sigma_k^{-1} \left[(x_i - \mu_k) (x_i - \mu_k)^T \right] \Sigma_k^{-1} \right] = 0$$

$$N_k \cdot I = \sum_{i=1}^n \gamma_t(z_i=k) \Sigma_k^{-1} \left[(x_i - \mu_k) (x_i - \mu_k)^T \right]$$

$$\Sigma_k^{+1} = \frac{1}{N_k} \sum_{i=1}^n \gamma_t(z_i=k) (x_i - \mu_k) (x_i - \mu_k)^T$$

$$p_4 \quad p(x | z=k) = \prod_{d=1}^D \theta_{kd}^{x_{kd}} (1 - \theta_{kd})^{1-x_{kd}}$$

$$p(z) = \prod_{k=1}^K \pi_k^{z_k} = \prod_{k=1}^K \left(\frac{1}{K}\right)^{z_k}$$

$$\textcircled{1} \quad \gamma(z_{i:k}) = \frac{\left(\frac{1}{K}\right)^{z_{i:k}} \cdot \theta_{jd}^{x_{jd}} (1 - \theta_{jd})^{1-x_{jd}}}{\sum_{j=1}^K \left(\frac{1}{K}\right)^{z_{j:k}} \cdot \theta_{jd}^{x_{jd}} (1 - \theta_{jd})^{1-x_{jd}}}$$

$$\textcircled{2} \quad E(\log p(x, z | \pi, \mu, \Sigma)) = \sum_{n=1}^N \sum_{k=1}^K \gamma_t(z_i=k) \ln p(x_i | z_i=k, \theta) \cdot p(z_i=k)$$

$$= \sum_{i=1}^n \sum_{k=1}^K \gamma_t(z_i=k) \left[\sum_{d=1}^D x_{id} \ln \theta_{kd} + (1 - x_{id}) \ln (1 - \theta_{kd}) + \ln \frac{1}{K} \right]$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = \sum_{i=1}^n \gamma_t(z_i=k) \left[x_{id} \frac{1}{\theta_{kd}} - (1 - x_{id}) \frac{1}{1 - \theta_{kd}} \right] = 0$$

$$\left[\frac{x_{id}}{\theta_{kd}} - \frac{1}{1 - \theta_{kd}} + \frac{x_{id}}{1 - \theta_{kd}} \right]$$

$$\left[\frac{x_{id} - x_{id} \theta_{kd} - \theta_{kd} + x_{id} \theta_{kd}}{\theta_{kd} (1 - \theta_{kd})} \right]$$

$$\sum_{i=1}^n \gamma_t(z_i=k) \frac{x_{id}}{\theta_{kd} (1 - \theta_{kd})} = \sum_{i=1}^n \gamma_t(z_i=k) \frac{\theta_{kd}}{\theta_{kd} (1 - \theta_{kd})}$$

$$\theta_{kd} = \frac{\sum_{i=1}^n \gamma_t(z_i=k) x_{id}}{\sum_{i=1}^n \gamma_t(z_i=k)}$$

Homework

Gaussian Mixture Model

Problem 5: Consider a mixture of K Gaussians

$$p(\mathbf{x}) = \sum_k \pi_k \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

Derive the expected value $\mathbb{E}[\mathbf{x}]$ and the covariance $\text{Cov}[\mathbf{x}]$.

Hint: it is helpful to remember the identity $\text{Cov}[\mathbf{x}] = \mathbb{E}[\mathbf{x}\mathbf{x}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^T$.

Problem 6: Consider a mixture of K isotropic Gaussians, all with the same *known* covariances $\boldsymbol{\Sigma}_k = \sigma^2 \mathbf{I}$.

Derive the EM algorithm for the case when $\sigma^2 \rightarrow 0$, and show that it's equivalent to Lloyd's algorithm for K -means.

Problem 7: Consider two random variables $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{y} \in \mathbb{R}^D$ distributed according to two different Gaussian mixture models with $\boldsymbol{\theta}^x = \{\boldsymbol{\pi}^x, \boldsymbol{\mu}^x, \boldsymbol{\Sigma}^x\}$ and $\boldsymbol{\theta}^y = \{\boldsymbol{\pi}^y, \boldsymbol{\mu}^y, \boldsymbol{\Sigma}^y\}$, i.e.

$$p(\mathbf{x} \mid \boldsymbol{\theta}^x) = \sum_{k=1}^{K_x} \pi_k^x \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k^x, \boldsymbol{\Sigma}_k^x),$$

$$p(\mathbf{y} \mid \boldsymbol{\theta}^y) = \sum_{l=1}^{K_y} \pi_l^y \mathcal{N}(\mathbf{y} \mid \boldsymbol{\mu}_l^y, \boldsymbol{\Sigma}_l^y),$$

and the random variable $\mathbf{z} = \mathbf{x} + \mathbf{y}$.

- Describe a generative process (process of drawing samples) for \mathbf{z} .
- Explain in a few sentences why $p(\mathbf{z} \mid \boldsymbol{\theta}^x, \boldsymbol{\theta}^y)$ is again a mixture of Gaussians.
- State the probability density function $p(\mathbf{z} \mid \boldsymbol{\theta}^x, \boldsymbol{\theta}^y)$ of \mathbf{z} .

Problem 8:

- Given is the dataset displayed in the figure below. Apply the K-means algorithm to this data using $K = 2$ and using the circled points as initial centroids.

What are the clusters after K-Means converges? Draw your solution in the figure above, i.e. mark the location of the centroids with \times 's and show the clusters by drawing two bounding boxes around the points assigned to each cluster.

How many iterations did it take for K-Means to converge in the above problem?

- Provide a different initialization, for which the algorithm will take **more** iterations to converge to the **same** solution. Make sure that your initialization does not lead to ties. Circle the initial centroids in the figure below.

$$P_5 \quad E(x) = E(E(x|z)) = \sum_{k=1}^K E(x|z) \cdot p(z) = \sum_{k=1}^K \mu_k \cdot \pi_k.$$

$$\begin{aligned} E(xx^T) &= E(E(xx^T|z)) = \sum_{k=1}^K E(xx^T|z) \cdot p(z) = \\ &= \sum_{k=1}^K (\text{cov}(x|z) + E(x|z) E(x|z)^T) \\ &= \sum_{k=1}^K (\Sigma_k + \mu_k \mu_k^T) \cdot \pi_k. \end{aligned}$$

$$\begin{aligned} \text{cov}(x) &= E(xx^T) - E(x) \cdot E(x)^T \\ &= \sum_{k=1}^K (\Sigma_k + \mu_k \mu_k^T) \pi_k - \sum_{k=1}^K \sum_{j=1}^K \mu_k \pi_k \mu_j^T \pi_j \end{aligned}$$

$$P_6. \quad p(z_k=1 | x_i, \theta) = \frac{1}{\sum_{k=1}^K \frac{\pi_k}{\pi_k} \exp\left(\frac{-\|x_i - \mu_k\|^2 + \|x_i - \mu_L\|^2}{2\sigma^2}\right)}$$

if μ_k is close to x_i , for $k=L$

$$\frac{-\|x_i - \mu_L\|^2 + \|x_i - \mu_k\|^2}{2\sigma^2} \leq 0$$

$k=L$ expl. = 1 $k \neq L$ exponentially increasing negative

if μ_k isn't close to x_i ($k \neq L$)

$$0 < \frac{-\|x_i - \mu_L\|^2 + \|x_i - \mu_k\|^2}{2\sigma^2} \rightarrow 0, \quad \sigma \rightarrow 0.$$

① Draw a sample x from $p(x|\theta^x)$ with the usual GMM

add them together

$$\textcircled{2} \quad p(z|\theta^x, \theta^y) = \sum_{k=1}^{K_x} \sum_{l=1}^{K_y} \pi_k^x \pi_l^y \mathcal{N}(z | \mu_k^x + \mu_l^y, \Sigma_k^x + \Sigma_l^y)$$

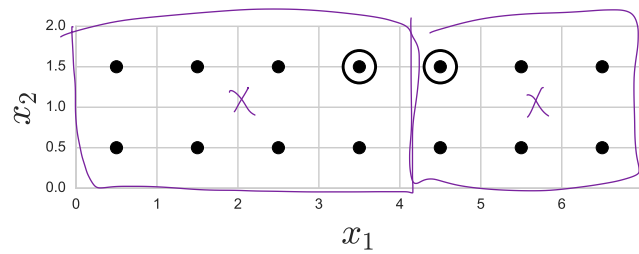


Figure 1: K-Means Dataset

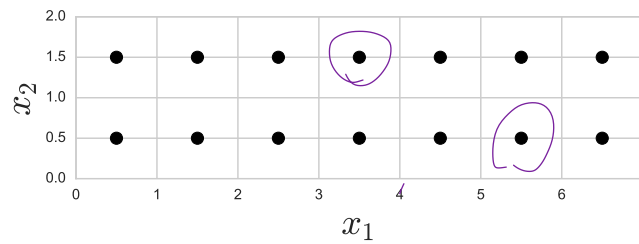


Figure 2: Provide your initialization

Problem 9: Download the notebook `exercise_12_clustering.ipynb` from Moodle. Fill in the missing code and run the notebook.