

# **Introduction to Machine Learning: Homework IV**

Due on Dec 7th, 2022 at 11:59pm

*Professor Ziping Zhao*

**Bingnan Li**  
2020533092

## 1. [Clustering and Mixture Models]

(a) K-means algorithm.

**Solution:**

- i. Initialize K cluster centers  $m_i$  by randomly selecting K input data points.
- ii. Repeat the following procedure until convergence:
  - A. For all  $x^{(l)} \in \mathcal{X}$ , we obtain the estimated labels

$$b_i^{(l)} = \begin{cases} 1, & \text{if } i = \arg \min_j \|x^{(l)} - m_j\| \\ 0, & \text{elsewhere} \end{cases}$$

- B. For all  $m_i$ , we obtain

$$m_i = \frac{\sum_l b_i^{(l)} x^{(l)}}{\sum_l b_i^{(l)}}$$

(b) Cluster the samples into 2 clusters.

**Solution:**

First, we select  $m_1 = (0, 0)$  and  $m_2 = (5, 0)$  as initialized cluster center. Then for the first iteration, we have the following result:

$$\begin{aligned} b_1^{(1)} &= 1 & b_2^{(1)} &= 0 \\ b_1^{(2)} &= 1 & b_2^{(2)} &= 0 \\ b_1^{(3)} &= 1 & b_2^{(3)} &= 0 \\ b_1^{(4)} &= 0 & b_2^{(4)} &= 1 \\ b_1^{(5)} &= 0 & b_2^{(5)} &= 1 \\ m_1 &= \frac{(0, 2) + (0, 0) + (1, 0)}{3} = \left(\frac{1}{3}, \frac{2}{3}\right) \\ m_2 &= \frac{(5, 0) + (5, 2)}{2} = (5, 1) \end{aligned}$$

Next, for the second iteration, we find that

$$\begin{aligned} b_1^{(1)} &= 1 & b_2^{(1)} &= 0 \\ b_1^{(2)} &= 1 & b_2^{(2)} &= 0 \\ b_1^{(3)} &= 1 & b_2^{(3)} &= 0 \\ b_1^{(4)} &= 0 & b_2^{(4)} &= 1 \\ b_1^{(5)} &= 0 & b_2^{(5)} &= 1 \\ m_1 &= \frac{(0, 2) + (0, 0) + (1, 0)}{3} = \left(\frac{1}{3}, \frac{2}{3}\right) \\ m_2 &= \frac{(5, 0) + (5, 2)}{2} = (5, 1) \end{aligned}$$

The result converged, so we terminated the algorithm and cluster centers are

$$m_1 = \left(\frac{1}{3}, \frac{2}{3}\right) \quad m_2 = (5, 1)$$

## 2. [Clustering and Mixture Models]

- (a) Advantages of GMM and Why it can be used for clustering.

**Solution:**

Advantages: GMM is a kind of "soft-label" method, the projected data do not represent deterministic classification label but the probability of belonging to any classes.

Why it can be used for clustering: K-means is a special case of GMM. In practice, the higher the  $h_i^{(l)}$  is, the more likely that  $x^{(l)}$  is generated by component  $\mathcal{G}_i$ , which can be interpreted as  $x^{(l)}$  belongs to cluster  $i$ .

- (b) Estimate the parameters of the GMM.

**Solution:**

Define  $\mathcal{Q}(\phi|\phi^t)$  as following

$$\begin{aligned}\mathcal{Q}(\phi|\phi^t) &= \mathbb{E}[\mathcal{L}_C(\phi|\mathcal{X}, \mathcal{Z})|\mathcal{X}, \phi^t] \\ &\text{where} \\ \mathcal{L}_C(\phi) &= \log \prod_l p(x^{(l)}, z^{(l)}|\phi) \\ &= \sum_l \left[ \log P(z^{(l)}|\phi) + \log p(x^{(l)}|z^{(l)}, \phi) \right] \\ &= \sum_l \sum_i z_i^{(l)} [\log \pi_i + \log p_i(x^{(l)}|\phi)]\end{aligned}$$

Hence

$$\begin{aligned}\mathcal{Q}(\phi|\phi^t) &= \mathbb{E}[\mathcal{L}_C(\phi|\mathcal{X}, \mathcal{Z})|\mathcal{X}, \phi^t] \\ &= \sum_l \sum_i \mathbb{E}[z_i^{(l)}|\mathcal{X}, \phi^t] [\log \pi_i + \log p_i(x^{(l)}|\phi)]\end{aligned}$$

where

$$\begin{aligned}\mathbb{E}[z_i^{(l)}|\mathcal{X}, \phi^t] &= \mathbb{E}[z_i^{(l)}|\mathbf{x}^{(l)}, \phi] \\ &= P(z_i^{(l)} = 1|\mathbf{x}^{(l)}, \phi^t) \\ &= \frac{p(\mathbf{x}^{(l)}|z_i^{(l)} = 1, \phi^t)P(z_i^{(l)} = 1|\phi^t)}{p(\mathbf{x}^{(l)}|\phi^t)} \\ &= \frac{p_i(\mathbf{x}^{(l)}|\phi^t)\pi_i}{\sum_j p_j(\mathbf{x}^{(l)}|\phi^t)\pi_j} \\ &= \frac{P(x^{(l)}|\mathcal{G}_i, \phi^t)\pi_i}{\sum_j P(x^{(l)}|\mathcal{G}_j, \phi^t)\pi_j} \\ &\equiv h_i^{(l)}\end{aligned}$$

Therefore, we have

$$\begin{aligned}h_i^{(l)} &= \frac{P(x^{(l)}|\mathcal{G}_i, \phi^t)\pi_i}{\sum_j P(x^{(l)}|\mathcal{G}_j, \phi^t)\pi_j} \\ &= \frac{|\Sigma_i|^{-\frac{1}{2}} \exp[-\frac{1}{2}(\mathbf{x}_l - \boldsymbol{\mu}_i)^T(\Sigma)^{-1}(\mathbf{x}_l - \boldsymbol{\mu}_i)]\pi_i}{\sum_{j=1}^K |\Sigma_j|^{-\frac{1}{2}} \exp[-\frac{1}{2}(\mathbf{x}_l - \boldsymbol{\mu}_j)^T(\Sigma)^{-1}(\mathbf{x}_l - \boldsymbol{\mu}_j)]\pi_j} \\ &= \frac{\mathcal{N}(\mathbf{x}_l|\boldsymbol{\mu}_i, \Sigma_i)\pi_i}{\sum_{j=1}^K \mathcal{N}(\mathbf{x}_l|\boldsymbol{\mu}_j, \Sigma_j)\pi_j}\end{aligned}$$

and

$$\mathcal{Q}(\phi|\phi^t) = \sum_l \sum_i h_i^{(l)} [\log \pi_i + \log \mathcal{N}(\mathbf{x}_l | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)]$$

Then, maximization of  $\mathcal{Q}(\phi|\phi^t)$  is equivalent to

$$\begin{aligned} & \underset{\{\pi_i\}, \{\boldsymbol{\mu}_i\}, \{\boldsymbol{\Sigma}_i\}}{\text{maximize}} && \mathcal{Q}(\phi|\phi^t) = \sum_l \sum_i h_i^{(l)} \log \pi_i + h_i^{(l)} \log \mathcal{N}(\mathbf{x}_l | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \\ & \text{subject to} && \sum_i \pi_i = 1 \end{aligned}$$

Since the second term does not depend on  $\pi_i$ , the problem for  $\{\pi_i\}$  is

$$\begin{aligned} & \underset{\{\pi_i\}}{\text{maximize}} && \sum_l \sum_i h_i^{(l)} \log \pi_i \\ & \text{subject to} && \sum_i \pi_i = 1 \end{aligned}$$

By using Lagrangian, we solve for

$$\frac{\partial}{\partial \pi_i} \left[ \sum_l \sum_i h_i^{(l)} \log \pi_i - \lambda \left( \sum_i \pi_i - 1 \right) \right] = 0$$

And we get

$$\pi_i = \frac{\sum_l h_i^{(l)}}{N}$$

Then the first term of  $\mathcal{Q}$  does not depend on  $\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i$ . Hence, the problem for  $\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i$  is

$$\underset{\{\boldsymbol{\mu}_i\}, \{\boldsymbol{\Sigma}_i\}}{\text{maximize}} \quad \sum_l \sum_i h_i^{(l)} \log \mathcal{N}(\mathbf{x}_l | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

By solving

$$\frac{\partial}{\partial \boldsymbol{\mu}_i} \left[ \sum_l \sum_i h_i^{(l)} \log \mathcal{N}(\mathbf{x}_l | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \right] = 0$$

and

$$\frac{\partial}{\partial \boldsymbol{\Sigma}_i} \left[ \sum_l \sum_i h_i^{(l)} \log \mathcal{N}(\mathbf{x}_l | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \right] = 0$$

we get

$$\boldsymbol{\mu}_i^{t+1} = \frac{\sum_l h_i^{(l)} \mathbf{x}_l}{\sum_l h_i^{(l)}}$$

and

$$\boldsymbol{\Sigma}_i^{t+1} = \frac{\sum_l h_i^{(l)} (\mathbf{x}_l - \boldsymbol{\mu}_i^{t+1})(\mathbf{x}_l - \boldsymbol{\mu}_i^{t+1})^T}{\sum_l h_i^{(l)}}$$

### 3. [Nonparametric Density Estimation]

(a) Expression of  $\hat{p}(x)$ .

**Proof:**

By definition, the histogram estimator is defined as following:

$$\hat{p}(x) = \frac{\#\{x^{(l)} \text{ in the same bin as } x\}}{nh} = \frac{\#\{x^{(l)} \in [\lfloor \frac{x}{h} \rfloor h, \lceil \frac{x}{h} \rceil h)\}}{nh}$$

- (b) Expression of
- $L'(h)$
- based on the histogram estimator
- $\hat{p}(x)$
- .

**Proof:**First, for the first term of  $L'$ , we can split the integral by bins:

$$\begin{aligned}\int_0^1 \hat{p}^2(x) dx &= \sum_{j=0}^{\lfloor \frac{1}{h} \rfloor} \int_{jh}^{(j+1)h} \frac{\#^2 \{x^{(l)} \in [jh, (j+1)h)\}}{n^2 h^2} dx \\ &= \frac{\sum_{j=0}^{\lfloor \frac{1}{h} \rfloor} \#^2 \{x^{(l)} \in [jh, (j+1)h)\}}{n^2 h}\end{aligned}$$

For the second term of  $L'$ , we can rewrite it as following:

$$\begin{aligned}\frac{2}{n} \sum_{i=1}^n \hat{p}(x_i) &= \frac{2}{n} \sum_{i=1}^n \frac{\# \{x^{(l)} \in [\lfloor \frac{x}{h} \rfloor h, \lceil \frac{x}{h} \rceil h)\}}{nh} \\ &= \frac{2 \sum_{j=0}^{\lfloor \frac{1}{h} \rfloor} \#^2 \{x^{(l)} \in [jh, (j+1)h)\}}{n^2 h}\end{aligned}$$

Hence, we have

$$L'(h) = \int_0^1 \hat{p}^2(x) dx - \frac{2}{n} \sum_{i=1}^n \hat{p}(x_i) = - \frac{\sum_{j=0}^{\lfloor \frac{1}{h} \rfloor} \#^2 \{x^{(l)} \in [jh, (j+1)h)\}}{n^2 h}$$

- (c)
- $h$
- that minimizes
- $L'(h)$
- .

**Proof:**

Since as  $h$  goes to 0, as long as there are not coincide sample points, the numerator of  $L'(h)$  will go to  $n$ . However, the denominator of  $L'(h)$  is  $n^2 h$ , which goes to 0 as  $h \rightarrow 0$ , thus  $\lim_{h \rightarrow 0} L'(h) = -\infty$ .

$$h = \arg \min_{h>0} L'(h) = 0$$

## 4. [Nonparametric Regression]

- (a) Estimated output
- $\hat{y}$
- and is linear regression a linear smoother?

**Solution:**Given that the least squares estimate for  $\mathbf{w}$  is

$$\mathbf{w}^* = (H^T H)^{-1} H^T Y$$

we have the following estimated output  $\hat{y}$ 

$$\begin{aligned}\hat{y} &= (\mathbf{w}^*)^T \cdot \mathbf{h}(\mathbf{x}) \\ &= \left[ (H^T H)^{-1} H^T Y \right]^T \mathbf{h}(\mathbf{x}) \\ &= Y^T H (H^T H)^{-1} \mathbf{h}(\mathbf{x}) \\ &= \left[ H (H^T H)^{-1} \mathbf{h}(\mathbf{x}) \right]^T Y \\ &= \mathbf{h}(\mathbf{x})^T (H^T H)^{-1} H^T Y \\ &\Rightarrow \\ \mathbf{l}(\mathbf{x}) &= H (H^T H)^{-1} \mathbf{h}(\mathbf{x})\end{aligned}$$

Hence, linear regression is a linear smoother.

- (b) In kernel regression, if we use kernel  $K(x_i, x) = \exp\left\{-\frac{\|x_i - x\|^2}{2\sigma^2}\right\}$ , given an input  $x$ , please derive the estimated output  $\hat{y}$ . Furthermore, is this kernel regression a linear smoother?

**Solution:**

By the definition of Kernel regression, we have the estimated output  $\hat{y}$  as following:

$$\begin{aligned}\hat{y} &= \frac{\sum_{i=1}^n K(x_i, x) y_i}{\sum_{i=1}^n K(x_i, x)} \\ &= \frac{\sum_{i=1}^n \exp\left\{-\frac{\|x_i - x\|^2}{2\sigma^2}\right\} y_i}{\sum_{i=1}^n \exp\left\{-\frac{\|x_i - x\|^2}{2\sigma^2}\right\}} \\ &= \sum_{i=1}^n \text{softmax}\left(\frac{\|x_i - x\|^2}{2\sigma^2}\right) y_i\end{aligned}$$

Then define that

$$S = \begin{bmatrix} \text{softmax}\left(\frac{\|x_1 - x\|^2}{2\sigma^2}\right) \\ \text{softmax}\left(\frac{\|x_2 - x\|^2}{2\sigma^2}\right) \\ \vdots \\ \text{softmax}\left(\frac{\|x_n - x\|^2}{2\sigma^2}\right) \end{bmatrix}$$

we have

$$\begin{aligned}\hat{y} &= \sum_{i=1}^n \text{softmax}\left(\frac{\|x_i - x\|^2}{2\sigma^2}\right) y_i \\ &= S^T Y \\ &\Rightarrow \\ l(x) = S &= \begin{bmatrix} \text{softmax}\left(\frac{\|x_1 - x\|^2}{2\sigma^2}\right) \\ \text{softmax}\left(\frac{\|x_2 - x\|^2}{2\sigma^2}\right) \\ \vdots \\ \text{softmax}\left(\frac{\|x_n - x\|^2}{2\sigma^2}\right) \end{bmatrix}\end{aligned}$$

Hence, this kernel regression is a linear smoother.