

CS280 Fall 2021 Assignment 1

Part A

ML Background

September 25, 2021

Name: 楊雨鑫
Student ID: 202023326

1. MLE (5 points)

Given a dataset $\mathcal{D} = \{x_1, \dots, x_n\}$. Let $p_{emp}(x)$ be the empirical distribution, i.e., $p_{emp}(x) = \frac{1}{n} \sum_{i=1}^n \delta(x, x_i)$ where $\delta(x, a)$ is the Dirac delta function¹ centered at a . Assume $q(x|\theta)$ be some probabilistic model.

- Show that $\arg \min_q KL(p_{emp}||q)$ is obtained by $q(x) = q(x; \hat{\theta})$, where $\hat{\theta}$ is the Maximum Likelihood Estimator and $KL(p||q) = \int p(x)(\log p(x) - \log q(x))dx$ is the KL divergence.

$$KL(p||q) = \int p(x)(\log p(x) - \log q(x))dx \\ = \int p(x) \log p(x) dx - \int p(x) \log q(x) dx$$

$$KL(p_{emp}(x)||q) = \int p_{emp}(x)(\log p_{emp}(x))dx - \int p_{emp}(x) q(x) dx$$

$$\arg \min_q KL(p_{emp}(x)||q) = \arg \min_q \left(\int p_{emp}(x) \log p_{emp}(x) dx \right. \\ \left. - \int p_{emp}(x) \log q(x) dx \right)$$

$$= \int p_{emp}(x) \log p_{emp}(x) dx - \operatorname{argmax}_q \int p_{emp}(x) \log q(x) dx$$

$$= \int p_{emp}(x) \log p_{emp}(x) dx - \int p_{emp}(x) \log \operatorname{argmax}_{\theta} q(x) dx$$

$\hat{\theta}$ is the maximum likelihood estimator.

$$\hat{\theta} = \operatorname{argmax}_{\theta} q(x)$$

So, $\arg \min_q KL(p_{emp}(x)||q)$ is obtained by $q(x) = q(x; \hat{\theta})$.

¹https://en.wikipedia.org/wiki/Dirac_delta_function

2. Gradient descent for fitting GMM (10 points)

Consider the Gaussian mixture model

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

where $\pi_j \geq 0$, $\sum_{j=1}^K \pi_j = 1$. (Assume $\mathbf{x}, \boldsymbol{\mu}_k \in \mathbb{R}^d, \boldsymbol{\Sigma}_k \in \mathbb{R}^{d \times d}$)

Define the log likelihood as

$$l(\theta) = \sum_{n=1}^N \log p(\mathbf{x}_n|\theta)$$

Denote the posterior responsibility that cluster k has for datapoint n as follows:

$$r_{nk} := p(z_n = k|\mathbf{x}_n, \theta) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k'} \pi_{k'} \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})}$$

- Show that the gradient of the log-likelihood wrt $\boldsymbol{\mu}_k$ is

$$\frac{d}{d\boldsymbol{\mu}_k} l(\theta) = \sum_n r_{nk} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

- Derive the gradient of the log-likelihood wrt π_k without considering any constraint on π_k .

(bonus 2 points: with constraint $\sum_k \pi_k = 1$)

$$\begin{aligned} \textcircled{1} \frac{d l(\theta)}{d P(\mathbf{x}|\theta)} &= \sum_{n=1}^N \frac{1}{P(\mathbf{x}_n|\theta)} \frac{d P(\mathbf{x}_n|\theta)}{d \boldsymbol{\mu}_k} = \frac{\sum_{k=1}^K \pi_k N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{d \boldsymbol{\mu}_k} = \sum_{k=1}^K \pi_k \frac{\partial N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ \frac{\partial N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} &= N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \end{aligned}$$

$$\frac{d l(\theta)}{d \boldsymbol{\mu}_k} = \sum_{n=1}^N \frac{1}{P(\mathbf{x}_n|\theta)} \frac{\partial P(\mathbf{x}_n|\theta)}{\partial \boldsymbol{\mu}_k} = \sum_{n=1}^N \frac{1}{P(\mathbf{x}_n|\theta)} \cdot \pi_k N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) = \sum_n r_{nk} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

$$\textcircled{2} \frac{d l(\theta)}{d \pi_k} = \sum_{n=1}^N \frac{1}{P(\mathbf{x}_n|\theta)} \frac{\partial P(\mathbf{x}_n|\theta)}{\partial \pi_k} = \sum_{n=1}^N \frac{N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k'} \pi_{k'} N(\mathbf{x}_n|\boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})}$$

when $\sum_k \pi_k = 1$, $\frac{\partial \pi_k}{\partial \pi_{k'}} = -1$ if $k \neq k'$. $\frac{\partial P(\mathbf{x}_n|\theta)}{\partial \pi_k} = N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) - \sum_{k' \neq k} N(\mathbf{x}_n|\boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})$

$$\frac{d l(\theta)}{d \pi_k} = \sum_{n=1}^N \frac{N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) - \sum_{k' \neq k} N(\mathbf{x}_n|\boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})}{\sum_{k'} \pi_{k'} N(\mathbf{x}_n|\boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})} \quad \text{when } \sum_k \pi_k = 1.$$