

# CS 446 / ECE 449 — Homework 4

*your NetID here*

Version 1.0

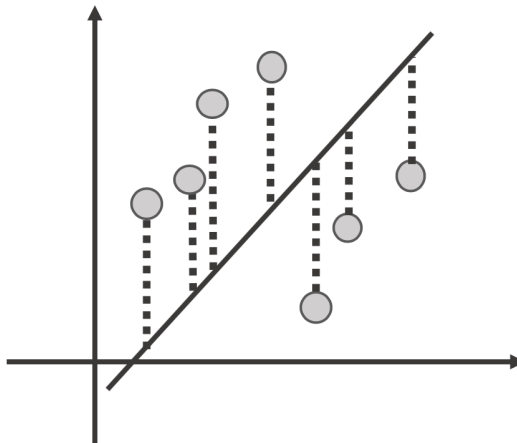
## Instructions.

- Homework is due **Tuesday, April 6th, at noon CST**; no late homework accepted.
- Everyone must submit individually at gradescope under **hw4** and **hw4code**.
- The “written” submission at **hw4** **must be typed**, and submitted in any format gradescope accepts (to be safe, submit a PDF). You may use L<sup>A</sup>T<sub>E</sub>X, markdown, google docs, MS word, whatever you like; but it must be typed!
- When submitting at **hw4**, gradescope will ask you to mark out boxes around each of your answers; please do this precisely!
- Please make sure your NetID is clear and large on the first page of the homework.
- Your solution **must** be written in your own words. Please see the course webpage for full academic integrity information. Briefly, you may have high-level discussions with at most 3 classmates, whose NetIDs you should place on the first page of your solutions, and you should cite any external reference you use; despite all this, your solution must be written in your own words.
- We reserve the right to reduce the auto-graded score for **hw4code** if we detect funny business (e.g., your solution lacks any algorithm and hard-codes answers you obtained from someone else, or simply via trial-and-error with the autograder).
- When submitting to **hw4code**, only upload **hw4.py** and **hw4\_utils.py**. Additional files will be ignored.

# 1. Principal Component Analysis

(a) For each of the following statements, specify whether the statement is true or false. If you think the statement is wrong, explain in 1 to 2 sentences why it is wrong.

- True or False: As shown in the figure below, PCA seeks a subspace such that the sum of all the vertical distance to the subspace (the dashed line) is minimized.



- True or False: PCA seeks a projection that best represents the data in a least-squares sense.
- True or False: PCA seeks a linear combination of variables such that the maximum variance is extracted from the variables.
- True or False: The principal components are not necessarily orthogonal to each other.
- True or False: Solving PCA using SVD might result in a solution which corresponds to a local minimum.

(b) Recall that PCA finds a direction  $w$  in which the projected data has highest variance by solving the following program:

$$\max_{w: ||w||^2=1} w^T \Sigma w. \quad (1)$$

Here,  $\Sigma$  is a covariance matrix. You are given a dataset of two 2-dimensional points  $(1, 3)$  and  $(4, 7)$ . Draw the two data points on the 2D plane. What is the first principal component  $w$  of this dataset?

(c) Now you are given a dataset of four points  $(2, 0)$ ,  $(2, 2)$ ,  $(6, 0)$  and  $(6, 2)$ . Draw the four data points on the 2D plane. Given this dataset, what is the dimension of the covariance matrix  $\Sigma$  in Eq. (1)? Also, explicitly write down the values of  $\Sigma$  given the dataset.

(d) What is the optimal  $w$  and the optimal value of the program in Eq. (1) given

$$\Sigma = \begin{bmatrix} 12 & 0 & 0 & 0 \\ 0 & 6 & 0 & 0 \\ 0 & 0 & 20 & 0 \\ 0 & 0 & 0 & 10 \end{bmatrix}$$

## 2. Gaussian Mixture Models

Consider a Gaussian mixture model with  $K$  components ( $k \in \{1, \dots, K\}$ ), each having mean  $\mu_k$ , variance  $\sigma_k^2$ , and mixture weight  $\pi_k$ . Further, we are given a dataset  $\mathcal{D} = \{x_i\}$ , where  $x_i \in \mathbb{R}$ . We use  $z_i = \{z_{ik}\}$  to denote the latent variables.

- (a) What is the log-likelihood of the data according to the Gaussian Mixture Model? (use  $\mu_k$ ,  $\sigma_k$ ,  $\pi_k$ ,  $K$ ,  $x_i$ , and  $\mathcal{D}$ ).
- (b) Assume  $K = 1$ , find the maximum likelihood estimate for the parameters  $(\mu_1, \sigma_1^2, \pi_1)$ .
- (c) What is the probability distribution on the latent variables, i.e., what is the distribution  $p(z_{i,1}, z_{i,2}, \dots, z_{i,K})$  underlying Gaussian mixture models. Also give its name.
- (d) For general  $K$ , what is the posterior probability  $p(z_{ik} = 1|x_i)$ ? To simplify, wherever possible, use  $\mathcal{N}(x_i|\mu_k, \sigma_k)$ , a Gaussian distribution over  $x_i \in \mathbb{R}$  having mean  $\mu_k$  and variance  $\sigma_k^2$ .
- (e) How are k-Means and Gaussian Mixture Model related? (There are three conditions)
- (f) Consider the modified Gaussian Mixture Model objective:

$$\min_{\mu} - \sum_{x_i \in \mathcal{D}} \epsilon \log \sum_{k=1}^K \exp(-(x_i - \mu_k)^2/\epsilon).$$

We denote by  $F_k = (x - \mu_k)^2$ . Show that:

$$\lim_{\epsilon \rightarrow 0} \epsilon \log \sum_{k=1}^K \exp(-F_k/\epsilon) = \sum_{k=1}^K \mathbb{1}_{\{k=\arg\min_k F_k\}} F_k, \quad \epsilon \in \mathbb{R}^+$$

**Hint:** Use l'Hopital rule.

- (g) Conclude that the objective for k-Means is the 0-temperature limit of Gaussian Mixture Model under the conditions you provided in part (e).

### 3. Expectation Maximization

In this problem, you will implement an expectation-maximization (EM) algorithm to cluster samples  $\mathcal{D} = \{x^{(i)}\}_{i=1}^n$ , with  $x^{(i)} \in \{0, 1\}^D$  into groups. You will be using a mixture of Bernoullis model to tackle this problem.

(a) **Mixture of Bernoullis.**

- i. Assume each variable  $x_d$  is drawn from a Bernoulli( $q_d$ ) distribution,  $P(x_d = 1) = q_d$ . Let  $q = [q_1, \dots, q_D] \in [0, 1]^D$  be the resulting vector of Bernoulli parameters. Write an expression for  $P(x|q)$  as a function of  $q_d$  and  $x_d$ .
- ii. Now suppose we have a mixture of  $K$  Bernoulli distributions: each vector  $x^{(i)}$  is drawn from some vector of Bernoulli random variables with parameters  $p^{(k)} = [p_1^{(k)}, \dots, p_D^{(k)}]$  that we call Bernoulli( $p^{(k)}$ ). Let  $\{p^{(1)}, \dots, p^{(K)}\} = p$ . Assume a distribution  $\pi$  over the selection of which set of Bernoulli parameters  $p^{(k)}$  is chosen. Write an expression for  $P(x^{(i)}|p, \pi)$ , as a function of  $\pi_k$  and  $P(x^{(i)}|p^{(k)})$ . Here  $\pi_k$  denotes the probability associated with the  $k^{\text{th}}$  Bernoulli component.
- iii. Using the above, write an expression for the log-likelihood of the data  $\mathcal{D}$ ,  $\log P(\mathcal{D}|\pi, p)$ .

(b) **Expectation step.**

- i. Let  $z^{(i)} \in \{0, 1\}^K$  be an indicator vector, such that  $z_k^{(i)} = 1$  if  $x^{(i)}$  was drawn from a Bernoulli( $p^{(k)}$ ), and 0 otherwise. Let  $Z = \{z^{(i)}\}_{i=1}^n$ . Write down the expression of  $P(z^{(i)}|\pi)$  as a function of  $\pi_k$  and  $z_k^{(i)}$ .
- ii. Write down the expression of  $P(x^{(i)}|z^{(i)}, p, \pi)$  as a function of  $P(x^{(i)}|p^{(k)})$  and  $z_k^{(i)}$ .
- iii. Using the two quantities above, derive the likelihood of the data and latent variables  $P(Z, \mathcal{D}|\pi, p)$ .
- iv. Let  $\eta(z_k^{(i)}) = \mathbb{E}[z_k^{(i)}|x^{(i)}, \pi, p]$ . Show that

$$\eta(z_k^{(i)}) = \frac{\pi_k \prod_{d=1}^D (p_d^{(k)})^{x_d^{(i)}} (1 - p_d^{(k)})^{1-x_d^{(i)}}}{\sum_j \pi_j \prod_{d=1}^D (p_d^{(j)})^{x_d^{(i)}} (1 - p_d^{(j)})^{1-x_d^{(i)}}}$$

- v. Let  $\tilde{p}$ ,  $\tilde{\pi}$  be the new parameters that we would like to maximize.  $p$ ,  $\pi$  are from the previous iteration. Use this to derive the following final expression for the E-step in the EM algorithm:

$$\mathbb{E}[\log P(Z, \mathcal{D}|\tilde{p}, \tilde{\pi})|\mathcal{D}, p, \pi] = \sum_{i=1}^n \sum_{k=1}^K \eta(z_k^{(i)}) \left[ \log \tilde{\pi}_k + \sum_{d=1}^D (x_d^{(i)} \log \tilde{p}_d^{(k)} + (1 - x_d^{(i)}) \log(1 - \tilde{p}_d^{(k)})) \right]$$

(c) **Maximization step.** In the following, we will find  $\tilde{p}$  and  $\tilde{\pi}$  that maximize the above expression.

- i. Show that  $\tilde{p}$  that maximizes the E-step is:

$$\tilde{p}^{(k)} = \frac{\sum_{i=1}^N \eta(z_k^{(i)}) x^{(i)}}{N_k},$$

where  $N_k = \sum_{i=1}^N \eta(z_k^{(i)})$ .

- ii. Prove that the value of  $\tilde{\pi}$  that maximizes the E-step is:

$$\tilde{\pi}_k = \frac{N_k}{\sum_{k'} N_{k'}}.$$

## 4. K-Means 1

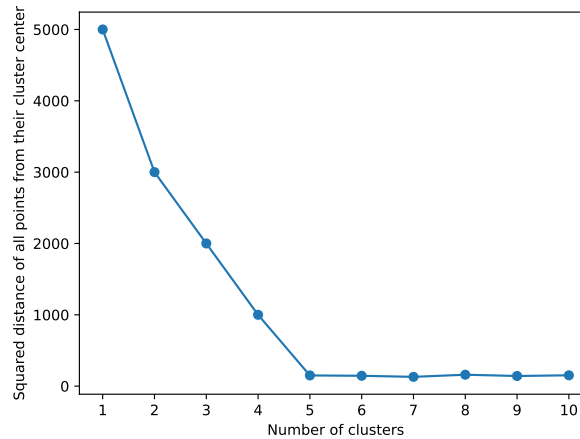
- (a) Mention if K-Means is a supervised or an un-supervised method and state the reason.
- (b) Assume that you are trying to cluster data points  $x_i$  for  $i \in \{1, 2, \dots, D\}$  into  $K$  clusters each with center  $\mu_k$  where  $k \in \{1, 2, \dots, K\}$ . The objective function for doing this clustering involves minimize the euclidean distance between the points and the cluster centers. It is given by

$$\min_{\mu} \min_r \sum_{i \in D} \sum_{k=1}^K \frac{1}{2} r_{ik} \|x_i - \mu_k\|_2^2.$$

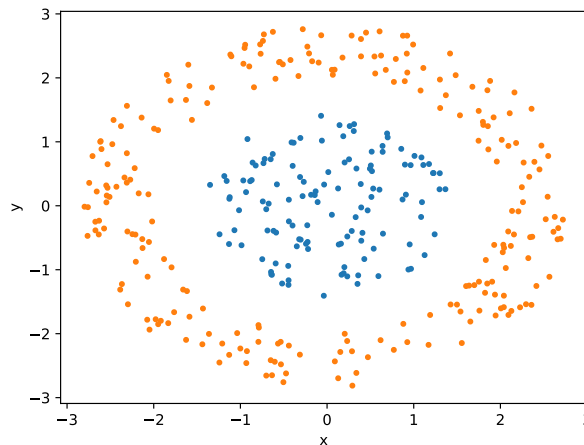
How do you ensure hard assignment of one data point to one and only one cluster at a given time?

**Hint:** By hard assignment we mean that you are 100 % sure that a point either belongs or doesn't belong to a cluster.

- (c) How does your answer to part b change if we want to obtain a soft assignment instead?  
**Hint:** By soft assignment we mean that a point belongs to a cluster with some probability.
- (d) Looking at the following plot, what is the best choice for the number of clusters?



- (e) Would K-Means be an efficient algorithm to cluster the following data? Explain your answer in a couple of lines.



## 5. K-Means 2

We are given a dataset  $\mathcal{D} = \{(x)\}$  of 2d points  $x \in \mathbb{R}^2$  which we are interested in partitioning into  $K$  clusters, each having a cluster center  $\mu_k$  ( $k \in \{1, \dots, K\}$ ) via the  $k$ -Means algorithm. This algorithm optimizes the following cost function:

$$\min_{\mu_k, r} \sum_{x \in \mathcal{D}, k \in \{1, \dots, K\}} \frac{1}{2} r_{x,k} \|x - \mu_k\|_2^2 \quad \text{s.t.} \quad \begin{cases} r_{x,k} \in \{0, 1\} & \forall x \in \mathcal{D}, k \in \{1, \dots, K\} \\ \sum_{k \in \{1, \dots, K\}} r_{x,k} = 1 & \forall x \in \mathcal{D} \end{cases} \quad (2)$$

- (a) What is the domain for  $\mu_k$ ?
- (b) Given fixed cluster centers  $\mu_k \forall k \in \{1, \dots, K\}$ , what is the optimal  $r_{x,k}$  for the program in Eq. 2? Provide a reason?
- (c) Given fixed  $r_{x,k} \forall x \in \mathcal{D}, k \in \{1, \dots, K\}$ , what are the optimal cluster centers  $\mu_k \forall k \in \{1, \dots, K\}$  for the program in Eq. 2?

**Hint:** Reason by first computing the derivative w.r.t  $\mu_k$ .

- (d) Using Pseudo-code, sketch the algorithm which alternates the aforementioned two steps. Is this algorithm guaranteed to converge and why? Is this algorithm guaranteed to find the global optimum? What is the reason?

**Hint:** you can provide a counter-example to invalidate a statement.

- (e) Please implement the aforementioned two steps. For the given dataset, after how many updates does the algorithm converge, what cost function value does it converge to and what are the obtained cluster centers? Visualize clusters at each step and attach the plots here. Please at least report numbers with one decimal point.

**Remark:** how we count updates: when computing a set of new centroids from initialization, we call this one update.

**Hint:** You may find `hw4.utils.vis_cluster` useful.