

CS 534 Machine Learning
Written Assignment 4
Aashish Adhikari

New Link

Date _____

Page _____

Q.1 Solution

We have our k-means objective

$$J = \sum_{c=1}^k \sum_{x_i \in C_c} |x_i - u_c|^2$$

Proving that it monotonically decreases means that we need to prove that it continuously decreases and eventually converges.

We know that if there are n data points, then we can have only k^n ways to make k clusters. During each new iteration, a new clustering is generated based on the original one. There are two possible cases for this: either the old and the new clustering are the same or the new clustering has a cost that is lower than the original. It is because each vector is assigned to the closest centroid, so the distance it contributes to the new cost decreases. Also, the cost decreases in the compilation step because the new centroid is a vector for which the cost reaches its minimum.

Since we have a finite set of possible clusterings, it will eventually arrive at a minimum.

We know that taking the partial derivative of the cost w.r.t the m^{th} components of the vector u

$$\frac{\partial J}{\partial u_m} = \sum_{x_i \in C_c} \partial(u_m - x_m)$$

Setting it to zero,

New Link

Date _____

Page _____

$$u_m = \frac{1}{|C_l|} \sum_{x_i \in C_l} x_m$$

Thus, we minimize the cost when the old centroid is replaced with the new centroid. Hence the sum of different components of the cost, J, must also decrease while computing the values.

$$Q.2 \quad J = \sum_{i=1}^k \sum_{x_i \in C_c} |x_i - u_c|^2$$

$$J' = \sum_{c=1}^{k+1} \sum_{x_i \in C_c} |x_i - u_c|^2$$

$$= \sum_{x_i \in C_{k+1}} |x_i - u_{k+1}|^2 + \sum_{c=1}^k \sum_{x_i \in C_k} |x_i - u_c|^2$$

In this new iteration, there is a point which belongs to C_{k+1} . For this particular point, the distance from itself to the centroid is zero. i.e., $|x_i - u_{k+1}| = 0$

In case more than one point are in this new cluster, the distance between the points and the new centroid will be less than the distance between the points and the centroid in the previous iteration k .

This implies

$$J = \sum_{c=1}^k \sum_{x_i \in C_c} |x_i - u_c|^2 < \sum_{c=1}^k \sum_{x_i \in C_c} |x_i - u_c|^2$$

Hence, the increase in the number of clusters, the resulting in a decrease in the objective J .

\Rightarrow Since the increase in the number of clusters will always reduce the distance to the datapoints, increasing k will always decrease this metric, that is, it reaches upto zero when its value is equal to the no. of data points. Hence, it is a BAD IDEA to choose the number of param clusters by minimizing J .

New Link

Date _____

Page _____

A solution to this is elbow method.

- Mean distance to the centroid as a function of k is plotted and the point, also called the elbow point, where the rate of decrease of sharply shifts can be roughly used to determine which k to use. An example of the plot is shown below:

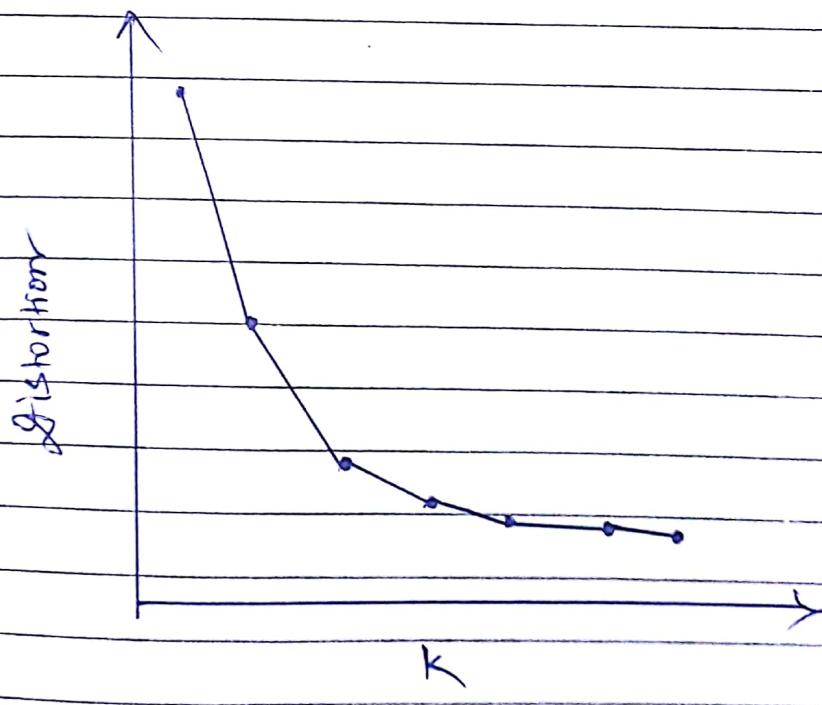


fig plot for distortion vs k

Q3

We have two univariate Gaussian distributions:

$$f(x|\theta_1) \text{ with } \mu_1 = \text{of } \sigma^2 = 1$$
$$f(x|\theta_2) \text{ with } \mu_2 = \text{of } \sigma^2 = 0.5$$

Then,

The mixture gives us

$$f(x) = \sum_{i=1}^k f(x, y=i)$$

$$= \sum_{i=1}^k f(x|y=i) f(y=i)$$

$$= \sum_{i=1}^k \alpha_i f(x|\theta_i)$$

In our case since we have only θ_1 & θ_2

$$= \alpha f(x|\theta_1) + (1-\alpha) f(x|\theta_2)$$

where α & $(1-\alpha)$ are class priors.

$$\theta_i = \{\mu_i, \Sigma_i\}$$

The log likelihood is

$$L(\theta; \mathbf{x}) = \sum_{i=1}^k \log [\alpha f(x|\theta_1) + (1-\alpha) f(x|\theta_2)]$$

direct maximization of $L(\theta | x)$ is difficult because of the sums of logarithmic terms.

So, we consider unobserved latent variables α_i :

$$\begin{aligned} & L(\theta; x, \Delta) \\ &= \sum_{i=1}^k \left[(1-\Delta_i) \log f(x|0_2) + \Delta_i \log f(x|\theta_i) \right] + \\ & \quad \sum_{i=1}^k \left[(1-\Delta_i) \log \alpha_i + \Delta_i \log (1-\alpha_i) \right] \end{aligned}$$

Since the value of Δ is unknown, we iterate the value of Δ with its expected value:

$$\gamma_i(\theta) = E(\Delta_i | \theta, x) = P(\Delta_i = 1 | \theta, x)$$

- Take initial guesses for parameters

$$m_1^1, \sigma_1^2; m_2^1, \sigma_2^2, \alpha$$

- Expectation step: Compute the responsibilities

$$\gamma_i^1 = \frac{\hat{\alpha} f(y_i)}{(1-\hat{\alpha}) f_0(y_i) + \hat{\alpha} f_1(y_i)}, \quad i=1, 2, \dots, k$$

$$\begin{aligned} & f_0^1(y_i) \\ &= f_{\hat{\alpha}_i}(x|\theta_i) \end{aligned}$$

- Maximization: compute the weighted mean & variances of the mixing probability

New Link

Date _____
Page _____

$$\Rightarrow \hat{\alpha} = \sum_{i=1}^k \frac{n_i}{K}$$

which is the required answer. #

Q4. We are given $p(x) = \sum_{k=1}^K \pi_k p(x|u_k)$

$$p(x|u_k) = \prod_{j=1}^M u_k(j)^{x(j)}$$

$$j = 1, \dots, M$$

$$x(j) = 1 \text{ if } x=j$$

E-step:

$$p(y_i=k|x_i) = \frac{p(x_i|y_i=k)p(y_i=k)}{p(x_i)}$$

$$= \frac{\pi_k \cdot p(x_i|u_k)}{\sum_{j=1}^K \pi_j p(x_i|u_j)}$$

$$= \frac{\pi_k \prod_{i=1}^M u_k(i)^{x(i)}}{\sum_{j=1}^K \pi_j \prod_{i=1}^M u_j(i)^{x_i(i)}}$$

$$\sum_{j=1}^K \pi_j \prod_{i=1}^M u_j(i)^{x_i(i)}$$

M step :

$$\arg \max_{\pi} \sum_{i=1}^N \sum_{j=1}^K p(y_i=j|x_i) \frac{\log p(x_i|y=j)}{p(y_i=j|x_i)}$$

$$= \sum_{i=1}^N \sum_{j=1}^K p(y_i=j|x_i) \log p(x_i|y_i=j) p(y_i=j)$$

$$= \sum_{i=1}^N \sum_{j=1}^K p(y_i=j|x_i) \log \prod_{k=1}^M \pi_j^{(k)} x_i^{(k)}$$

$$= \sum_{i=1}^N \sum_{j=1}^K p(y_i=j|x_i) (\log \pi_j + \sum_{k=1}^M \pi_i^{(k)} \log \pi_k)$$

This is our first equation (A)

Taking only π_T terms:

$$\sum_{i=1}^N p(y_i=t|x_i) \log \pi_t$$

Now we can use Lagrangian with a constraint

$$\sum_{j=1}^K \pi_j = 1$$

$$Z(\pi_T) = \sum_{i=1}^N p(y_i=t|x_i) \log \pi_t + \beta \left(\sum_{j=1}^K \pi_j - 1 \right)$$

$$\frac{\partial Z(\pi_T)}{\partial \pi_T} = \sum_{i=1}^N p(y_i=t|x_i) \cdot \frac{1}{\pi_t} + \beta = 0$$

$$\pi_T = -\frac{1}{\beta} \sum_{i=1}^N p(y_i=t|x_i) \quad (B)$$

we know that

$$\sum_{j=1}^K \pi_j = 1$$

$$\sum_{j=1}^K \left(-\frac{1}{\beta} \sum_{i=1}^N p(y_i=j|x_i) \right) = 1$$

$$-\beta = \sum_{j=1}^K \sum_{i=1}^N P(y_i=j | x_i) = N$$

Plugging the values back in \textcircled{B}

$$\pi_t = \frac{\sum_{i=1}^N P(y_i=t | x_i)}{N} \quad \textcircled{C}$$

Taking the terms of μ_t from \textcircled{A} , we get

$$\sum_{i=1}^N \left\{ P(y_i=t_k) \sum_{k=1}^M \pi_i(k) \log \mu_t(k) \right\}$$

Using the Lagrangian

$$\sum_{j=1}^M \mu_t(j) = 1$$

$$\begin{aligned} L(\mu_t) = & \sum_{i=1}^N \left\{ P(y_i=t_k) \sum_{k=1}^M \pi_i(k) \log \mu_t(k) \right. \\ & + \left. P \left(\sum_{j=1}^M \mu_t(j) - 1 \right) \right\} \end{aligned}$$

Taking derivative wrt $\mu_t(k)$

$$\frac{\partial}{\partial \mu_t(k)} = \sum_{i=1}^N \left(P(y_i=t | x_i) \cdot \pi_i(k) \frac{1}{\mu_t(k)} \right) + B = 0$$

$$\mu_t(k) = \left[\sum_{i=1}^N p(y_i = t | x_i) \cdot x_i(k) \right] \frac{1}{-\beta} \quad \textcircled{A}$$

We know $\sum_{j=1}^M \mu_t(j) = 1$

This gives,

$$\sum_{j=1}^M \sum_{i=1}^N p(y_i = t | x_i) x_i(j) \frac{1}{-\beta} = 1$$

or, $-\beta = \sum_{j=1}^M \sum_{i=1}^N p(y_i = t | x_i) x_i(j)$

From \textcircled{A}

$$\mu_t(k) = \frac{\sum_{i=1}^N p(y_i = t | x_i) x_i(k)}{\sum_{j=1}^M \sum_{i=1}^N p(y_i = t | x_i) x_i(j)}$$

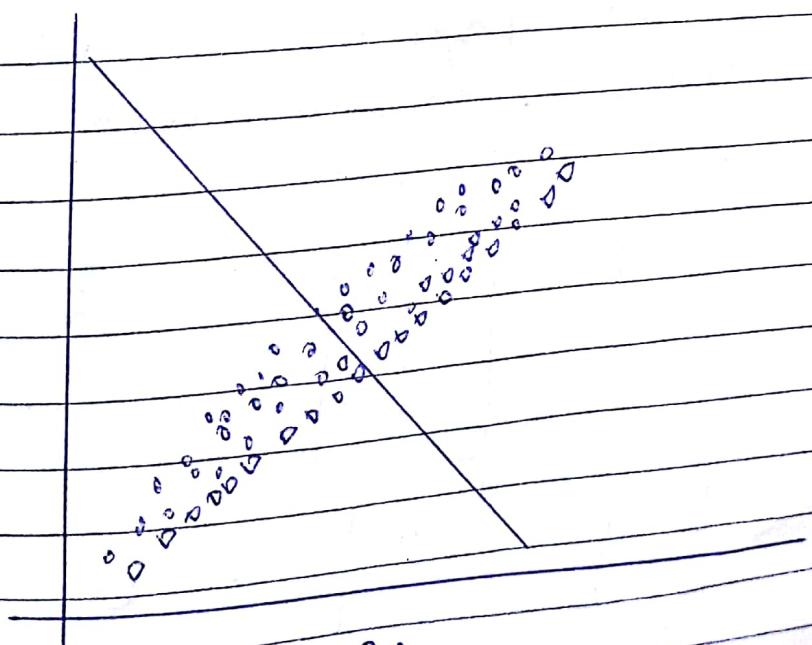
$$\sum_{j=1}^M \sum_{i=1}^N p(y_i = t | x_i) x_i(j) = 1$$

However, $\sum_{j=1}^M x_i(j) = 1$

$$\Rightarrow \mu_t(k) = \frac{\sum_{i=1}^N p(y_i = t | x_i) x_i(k)}{\sum_{i=1}^N p(y_i = t | x_i)} \quad \text{if}$$

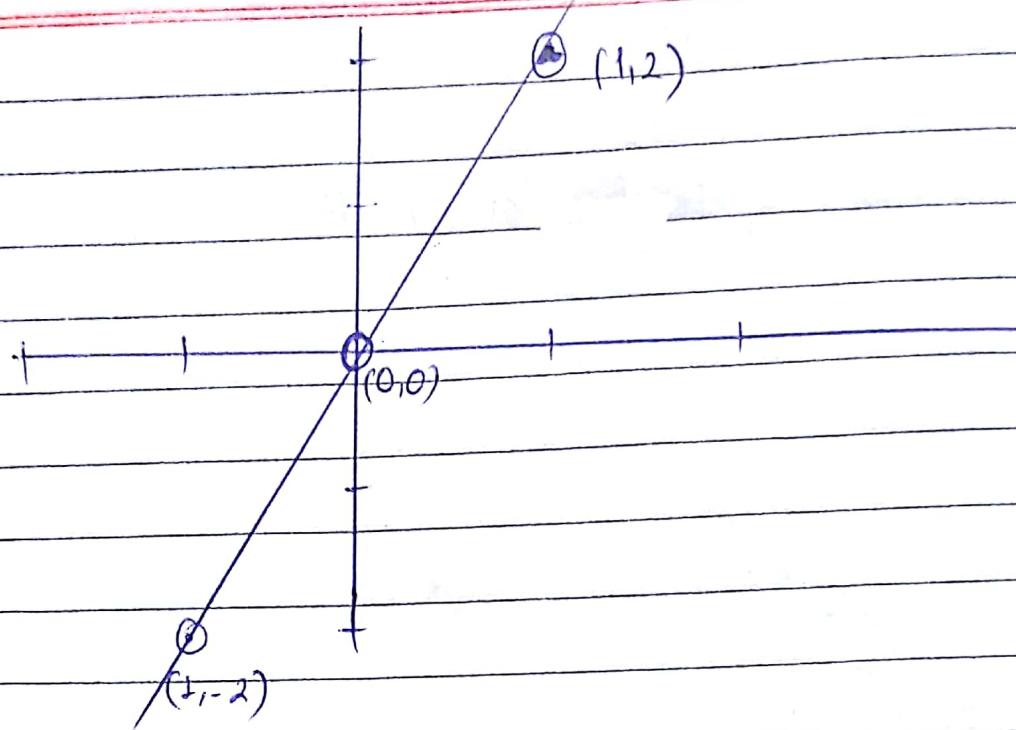


First Principal Component Direction



LDA

b)



Our data points are $(0,0), (1,2), (-1,-2)$.

Equation of the line is $y - 2x = 0$

First Principal Component direction would be along this line $y - 2x = 0$

$$w = \begin{pmatrix} 1/\sqrt{5} & 2/\sqrt{5} \end{pmatrix}^T$$

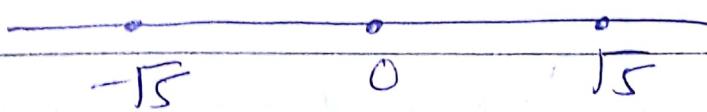
This vector is used to project the data points in a new one-dimensional space as $w^T x$.

$$\Rightarrow \begin{bmatrix} 1/\sqrt{5} & 2/\sqrt{5} \end{bmatrix} \begin{bmatrix} -1 \\ -2 \end{bmatrix} = -\sqrt{5}$$

$$\Rightarrow \begin{bmatrix} 1/\sqrt{5} & 2/\sqrt{5} \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} = 0$$

$$\Rightarrow \begin{bmatrix} 1/\sqrt{3} & 2/\sqrt{3} \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \sqrt{5}$$

In 1-D space



Variance of the projected data $= \frac{1}{N} \sum (x - \mu)^2$

$$\mu = (\sqrt{5} + 0 + (-\sqrt{5})) / 3 \\ = 0$$

$$\therefore \sigma^2 = \frac{1}{3} \times \sum (x - \mu)^2$$

$$= \frac{1}{3} [(\sqrt{5})^2 + 0 + (-\sqrt{5})^2]$$

$$= 10/3 \quad \#$$