# Chapter 12 Exercises

I Recitation Exercises

(1)a) We know that:

$$\frac{1}{|C_k|} \sum_{i,i' \in C_k} \sum_{j=1}^{P} (x_{ij} - x_{i'j})^2 = 2 \sum_{i \in C_k} \sum_{j=1}^{P} (x_{ij} - \bar{x}_{kj})^2$$

Let's decompose the left part

$$\sum_{i,i' \in C_k} \sum_{j=1}^{P} (x_{ij}^2 - 2x_{ij} + x_{i'j}^2) = |C_k| \sum_{i \in C_k} \sum_{j=1}^{P} x_{ij}^2 - 2 \sum_{i \in C_k} \sum_{i' \in C_k} \sum_{j=1}^{P} x_{ij} x_{i'j}$$

$$+ |C_k| \sum_{i \in C_k} \sum_{j=1}^{P} x_{ij}^2$$

Let us focus on middle term

$$-2 \sum_{i \in C_k} \sum_{i' \in C_k} \sum_{j=1}^{P} x_{ij} x_{i'j} = -2 \sum_{j=1}^{P} \sum_{i \in C_k} x_{ij} \left( \sum_{i' \in C_k} x_{i'j} \right)$$

But we also know that

$$C_k : \bar{x}_{kj} = \frac{1}{|C_k|} \sum_{i \in C_k} x_{ij}$$

So, finally we the part as:

$$2|C_k| \sum_{i \in C_k} \sum_{j=1}^{P} x_{ij}^2 - 2|C_k|^2 \sum_{j=1}^{P} \bar{x}_{kj}^2$$

Now divide $|C_k|$ and

$$2 \sum_{i \in C_k} \sum_{j=1}^{p} x_{ij}^2 - 2|C_k| \sum_{j=1}^{p} \bar{x}_{kj}^2$$

Rewrite the term as

$$-2|C_k| \sum_{j=1}^{p} \bar{x}_{kj}^2 = -4|C_k| \sum_{j=1}^{p} \bar{x}_{kj}^2 + 2|C_k| \sum_{j=1}^{p} \bar{x}_{kj}^2 = -4 \sum_{i \in C_k} \sum_{j=1}^{p} x_{ij} \bar{x}_{kj}$$
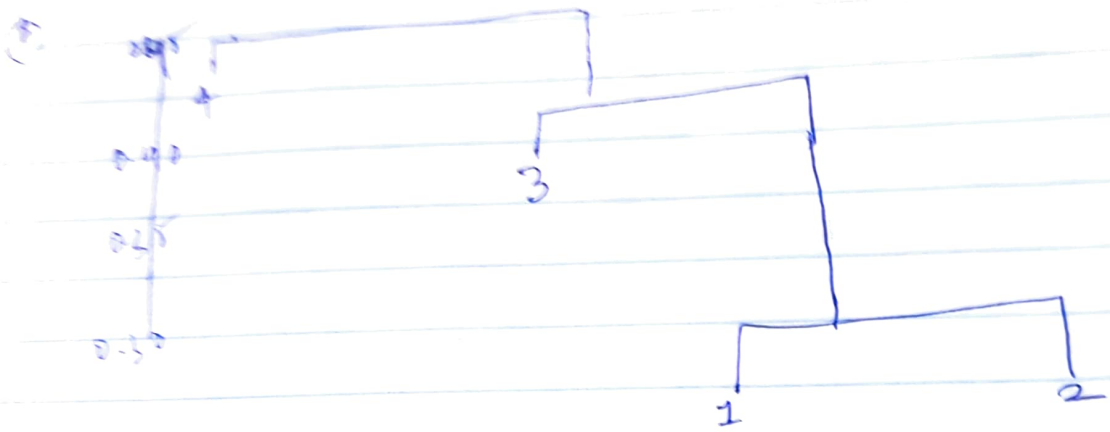
$$+ 2 \sum_{i \in C_k} \sum_{j=1}^{p} \bar{x}_{kj}^2$$

Replace with the rewritten left term and simplify

$$2 \sum_{i \in C_k} \sum_{j=1}^{p} x_{ij}^2 - 4 \sum_{i \in C_k} \sum_{j=1}^{p} x_{ij} \bar{x}_{kj} + 2 \sum_{i \in C_k} \sum_{j=1}^{p} \bar{x}_{kj}$$
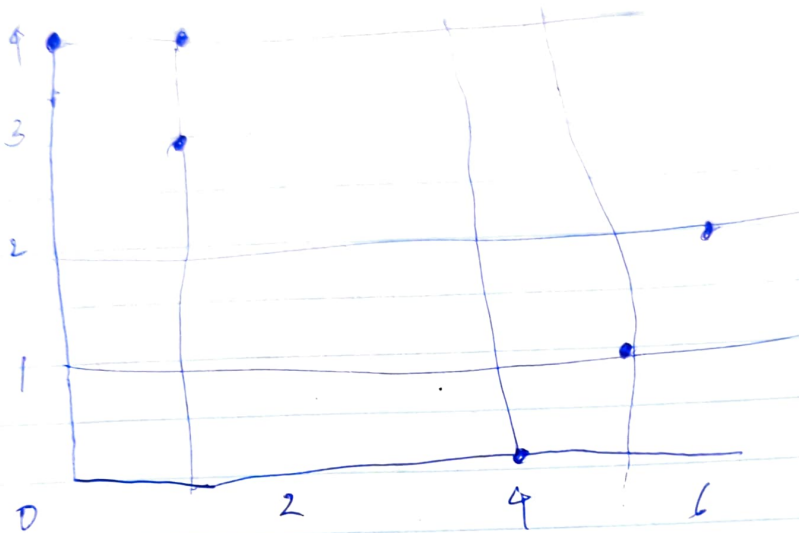
$$= 2 \sum_{i \in C_k} \sum_{j=1}^{p} (x_{ij} - \bar{x}_{kj})^2$$

(b) In k means clustering minimization of sum of squared euclidean distance for each cluster and our objective is the minimize the distance with cluster variance for each cluster and both cove same.

(i)

$0.90$

$0.75$

$0.60$

$0.50$

3

1    2

① Obcservations 1 and 2 will be in same cluster and observation

cluster diagram



(b) The random clusters are

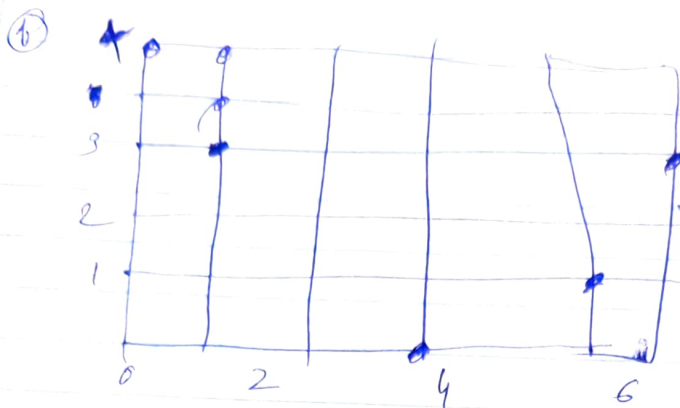$$[1 \ 1 \ 1 \ 1 \ 2 \ 2]$$

(d) After the iterations the 3 points near the left corner are assigned to same, on right corner are another clust

$$[1, 1, 1, 2, 2, 2]$$

(c) Centroids of each cluster

| Cluster | X1 | X2 |
|---------|------|----|
| 1 | 1.75 | 3 |
| 2 | 5.00 | 1 |

© At iteration 2, it is the same cluster

ⓑ


(Exercise 4)

ⓐ In single linkage clustering, the distance between two clusters is defined as the minimum distance between any two observations belonging to different cluster, in complete linkage clustering is defined as maximum & minimum values on equal, will fuse at same height, else the single linkage will fuse at a lower height. But without knowing this information we cant concluded so we can only say we need more information.

ⓑ Since we are fusing singleton cluster in both single linkage and complete linkage, the minimum distance between observation is equal to maximum distance

- Same height in both dendograms.