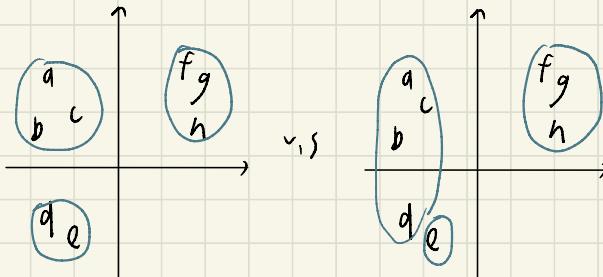
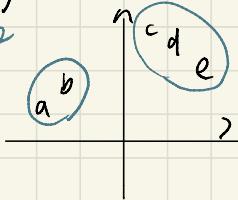


Clustering

$n=5$

$p=2$

把 data 分群



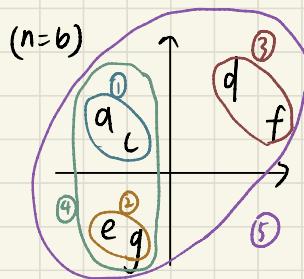
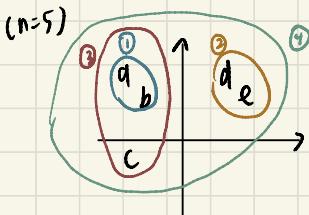
不同的分法？

Goal of clustering

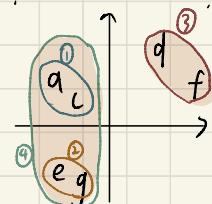
- ① Group similar observations
組內差異小，組間差異大
- ② Ideally, each group is very different.

How to generalize to $p > 2$?

1. Hierarchical clustering



想分 2 組 不要取第 5 步
3 : 4 :

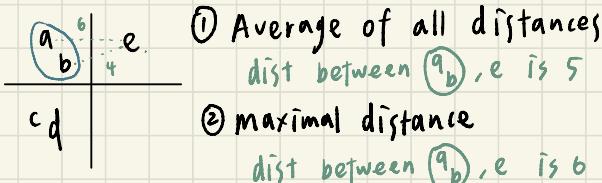


- Algorithm:

Bottom up

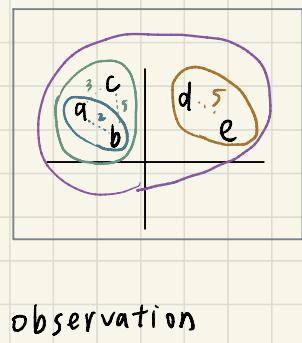
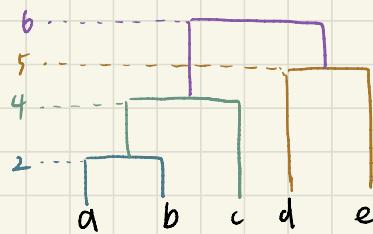
- repeat
- ① Calculate pairwise distance
 - ② Fuse the two nearest points
 - ③ repeat ①, ② until every observation is in one group.

- How to calculate distance between groups?



- Dendrogram

y = distance
between points
fused



2. K-means clustering

- Goal = find clusterings C_1, C_2, \dots, C_K (K clusters)

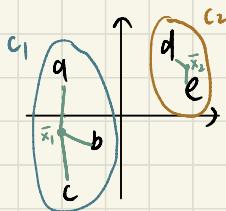
C_1, C_2, \dots, C_K is a partition of $\{1, 2, \dots, n\}$ (e.g. $C_1 = \{1, 3, 5\}$)

- $C_K \cap C_{K'} = \emptyset$
 - $C_1 \cup C_2 \cup \dots \cup C_K = \{1, 2, \dots, n\}$
- every observation is assigned and only assigned to one group

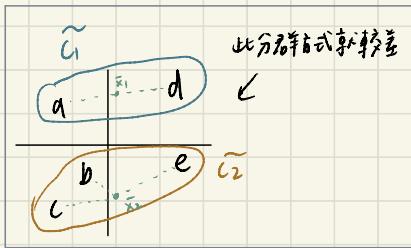
- K-means clustering find optimal C_1, C_2, \dots, C_K

$$\min \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - \bar{x}_k\|^2$$

, where $\bar{x}_k = \frac{1}{n_k} \sum_{x_i \in C_k} x_i$ (centroids)
 $n_k = \# \text{observation in } C_k$



Find clustering s.t. within group variation is the smallest.



Solution space = all possible partition C_1, C_2, \dots, C_K

Ques: How many possibilities?

- ① n^K n 等於 K 個選擇
- ② K^n n ; every observation has K choices
- ③ nK NP-Hard.

Hence, This is a combinatorial problem (NP Hard)

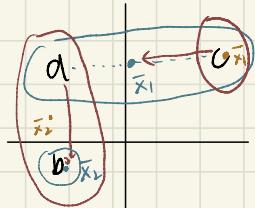
Remark

optimization is general hard outside calculus/linear programs

• K-means Algorithm

- ① Randomly assign points to cluster
- ② Calculate centroid $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_K$
- ③ Re-assign points to nearest centroids.) repeat until convergence.

$n=3, k=2$



- ① random assign, calculate x_i
- ② reassign
- ③ calculate x_i

Remark :

- ① it always converges
- ② But does not guarantee to the optimum
So, try many time with different initialization