① K-means.　clustering
② hierarchical clustering　( computational intensive
　　　　　　　　　　　　　　　　for large datasets)

NOTICE¡ Data Normalization maybe Required if attributes ranges
　　　　　　　　　very significantly.

⚡ Steps for K-means:

　　① define # of K.
　　② randomly assign obs to cluster centers $C_1, C_2, \cdots C_k$.
　　③ calculate the centroids of cluster centers $\mu_1, \mu_2, \cdots \mu_k$.
　　　　　　(mean)
　　④ calculate distance between each observation $x_1, x_2, \cdots x_n$ and
　　　　　　　　　　　　　　　cluster centroids $\mu_1, \mu_2, \cdots \mu_k$.

　　　　$d(x_i, \mu_i) = (x_1 - x_2)^2 + (y_1 - y_2)^2 + \cdots$　　　or $\sqrt{\cdots}$

(sum of squared error)⟩　　　　　　　　　(to reduce computation without sqrt)
　　⑤ calculate $SSE_T = SSE_1 + SSE_2 + \cdots + SSE_k$

　　⑥ iterate (repeat) above until $SSE_T$ dose not decrease.

　　　　(choosing randomly initial points (centroids) are crucial for K-means)

⚡ Steps for Hierarchical clustering.　　　　　⟋squared euclidian (maybe this
　　　　　　　　　　　　　　　　　　　　　⟋〈euclidian.　　　　in exam)
　　① calculate point-wise distances. (Euclidan / Manhattan)
　　　　　　　　　　　　　　⤷ generate proximity Matrix.
　　② fuse (merge) the closet. two.

　　③ calculate distance again. (either point-to-point ← euclidant/manhattu
　　　　　　and update matrix.　　or　　point-to-cluster ← linkage
　　　　　　　　　　　　　　　　　　　cluster-to-cluster ← linkage
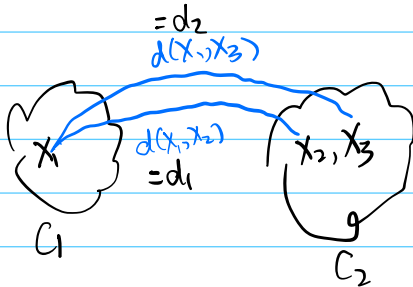　④ fuse the closet two.

⑤ repeat until we reach the root of dendrogram.

how to calculate distances
$$\begin{cases} \text{squared euc} : (x_a - x_b)^2 + (y_a - y_b)^2 \\ \text{euc} : \sqrt{(x_a - x_b)^2 + (y_a - y_b)^2} \\ \text{manhattan} : |x_a - x_b| + |y_a - y_b| \end{cases}$$

• $(x_a, y_a)$ ———— • $(x_b, y_b)$

how to calculate point/cluster - to - cluster linkage ? :



$$d(C_1, C_2) = \begin{cases} \text{single} : \min(d_1, d_2) \\ \text{complete} : \max(d_1, d_2) \\ \text{average} : \dfrac{d_1 + d_2}{2} \\ \text{centroid} : d(\mu_1, \mu_2) \end{cases}$$

(Single)
min $\begin{cases} \triangle : \text{can handle non-elliptical shapes} \\ \triangledown : \text{sensitive to noise and outliers} \end{cases}$

(Complete)
max $\begin{cases} \triangle : \text{less sensitive to noise and outliers ①} \\ \triangledown : \begin{cases} \text{tends to break large clusters} \\ \text{biased towards globular clusters ②} \end{cases} \end{cases}$

average $\begin{cases} \triangle : \text{same as ① (noise-outliers)} \\ \triangledown : \text{same as ② ( biased )} \end{cases}$

k-means vs hierarchical ?

k-means $\begin{cases} \triangle \text{ fast .} \\ \triangledown \text{ because it uses } \underline{mean} \text{ . easy to be affected} \\ \qquad \qquad \qquad \qquad \qquad \underline{\text{by the outliers}} \end{cases}$

✩ How we determind the best cluster?

(by using cluster index — silhouette index)



(a-value)

$a_1$ = average distance from point 1 to other points at the same cluster.

$$a_1 = avg.(\ d(x_1,x_2),\ d(x_1,x_3),\ d(x_1,x_4))$$

(b-value)

$b_1$ = minimum of average distance from point 1 to other clussters

① calc avg distance from $x_1$ to other clusters.
② took the minimum.

$$d(x_1,②) = \frac{d(x_1,y_1)+d(x_1,y_2)+d(x_1,y_3)}{3}$$

$$d(x_1,③) = \frac{d(x_1,z_1)+d(x_1,z_2)+d(x_1,z_3)}{3}$$

$$b_1 = min\ (\ d(x_1,②),\ d(x_1,③))$$

(S-value)

$$S_i = \frac{(b_i - a_i)}{\max(b_i, a_i)} \cdot \text{ generally } b_i > a_i$$

$$\text{ii} \longrightarrow S_i = \frac{b_i - a_i}{b_i}$$

$\dot{\rightarrow}f$ within cluster distance is very small $\longrightarrow a_i \rightsquigarrow 0$
then $S_i \rightsquigarrow 1$. (very good!)

$\dot{v}f$ within cluster distance is almost same as inter-cluster $\longrightarrow a_i \approx b_i$
then $S_i \longrightarrow 0$. (very bad, means that data point
actually belongs to another cluster)

$\dot{v}f$ a value larger than $b_i$ then $\max(b_i, a_i) \longrightarrow a_i$
in extreme $b_i - a_i \rightsquigarrow -a_i$.
ii $S_i \longrightarrow -1$.

ii in general i $-1 < \text{Svalue} \leq 1$

① we calculate S-value for each point

② we calculate S-value for each cluster $S_k = \frac{1}{n_k} \sum S_i$ $\quad$ n item in k.

③ we calculate S-value for the clustering model $I = \frac{1}{k} \sum S_k$

(sihoutte Index) $\qquad$ total # of clusters

close to '1 is better.