

Cluster Analysis

Grouping of objects (with multidimensional feature vectors) based on self similarities ("closeness").

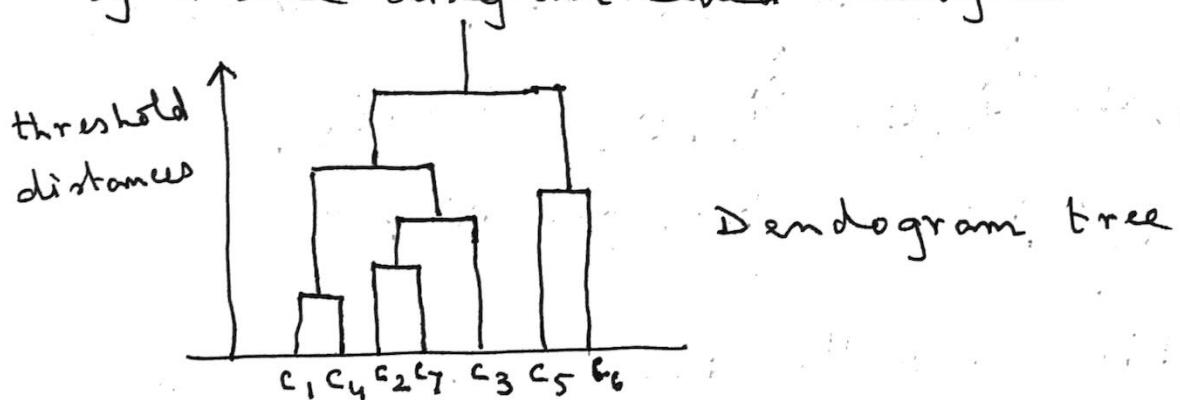
Two approaches

- Hierarchical clustering
- non-hierarchical clustering

Hierarchical clustering

In HCA the observation vectors (cases) are grouped together on the basis of mutual distances in such a manner so as to ensure hierarchy in cluster formation.

HCA output is usually visualized through a hierarchical tree, which is a nested set of partitions represented by a tree diagram called dendrogram.



Characteristics

- (i) Sectioning a tree at a particular level produces a partition into g disjoint groups
- (ii) If two groups are chosen from different partitions (i.e. the results of partitioning at different levels) then either the groups are disjoint or one group contains the other

- (iii) a numerical value is associated with each partition up/down the tree where branches join. This is a measure of distance between 2 merged clusters.
- (iv) We get different trees corresponding to different distance measures ("bet" cases).

Two algorithms for HCA

- Agglomerative clustering algorithm
- Divisive clustering algorithm

Agglomerative clustering algorithm

- Operates with successive merger of objects
- Starts with n clusters, each containing a single data point.
- At each stage merges the 2 most similar groups to form a new cluster, thus reducing the number of clusters by 1.
- The process of merger of group of objects continue till all subgroups are fused together to form a single cluster

Divisive clustering algorithm

- Operates by successively splitting groups of objects
- Starts with a single cluster having all the n objects
- Initial group split into 2 clusters \Rightarrow the distance bet["] the split groups ~~are~~ is maximum, i.e. the subgroups are as far as possible.

- Process of splitting continues till there are n clusters, each having a single object (leaf node)
- DCA is computationally inefficient.

Note: Result from agglomerative (or divisive) clustering algorithm displayed through dendrogram.

Steps in agglomerative hierarchical algorithm

Step I : n clusters, each having a single object and an $n \times n$ symmetric matrix of distances ($D = ((d_{ij}))$)

Step II : Search the distance matrix for the nearest (most similar) pair of clusters
Let the distance betⁿ the most similar clusters, say U and V , be d_{UV} .

Step III : Merge $U \& V$ to form cluster (U, V) , at a merger (fusion) level d_{UV} .

Update the distance matrix by

(i) deleting the rows & col^ms corresponding to clusters U & V

and (ii) adding a row & col^m giving the distance betⁿ (U, V) and the remaining clusters (existing ones other than $U \& V$)..

Step IV : Repeat (II) & (III) $n-1$ times, recording the identity of the clusters that are merged and the merger levels (i.e. the

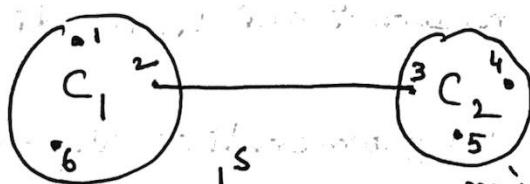
distances at which mergers take place).

Step 2 : Construct the dendrogram tree from the information on mergers and level of mergers.

Remark: step iii requires calculation of distance bet" clusters (for distance matrix updation).

The following are commonly used distance measures bet" clusters of objects

(i) Single link distance (minimum distance or nearest neighbor distance)



$$d_{c_1, c_2} = \min_{\substack{i \in c_1 \\ j \in c_2}} d_{ij}$$

(ii) Complete link distance (maximum distance or farthest) distance

$$d^c_{c_1, c_2} = \max_{\substack{i \in c_1 \\ j \in c_2}} d_{ij}$$

(iii) Average link distance

N_{c_1} : # of objects in c_1

N_{c_2} : # of objects in c_2

$$d^A_{c_1, c_2} = \frac{1}{N_{c_1} N_{c_2}} \sum_{\substack{i \in c_1 \\ j \in c_2}} d_{ij}$$

(iv) Centroid distance : Distance betⁿ cluster means

(v) Median distance : Distance betⁿ cluster medians

Remark : Single linkage agglomerative algorithm

In step iii (distance matrix update stage), we compute the distance betⁿ (U, V) and any existing cluster W as

$$d_{(U,V),W} = \min(d_{UW}, d_{VW})$$

↑ ↑
present in prev round

d_{UW} & d_{VW} are nearest neighbor distances

Remark : Complete linkage agglomerative algorithm

In step iii, compute distance betⁿ (U, V) and any existing W as

$$d_{(U,W),W} = \max(d_{UW}, d_{VW})$$

Remark : Average linkage agglomerative algorithm

In step iii, compute distance betⁿ (U, V) and W computed as

$$d_{(U,V),W} = \frac{\sum_{i \in (U,V)} \sum_{k \in W} d_{ik}}{N_{(U,V)} N_W}$$

d_{ik} : distance betⁿ object i in the cluster (U, V)
and object k in cluster W

$N_{(U,V)}$: # of objects in (U, V) & N_W : # of objects in W

Example: Single linkage clustering

Five objects c_1, c_2, c_3, c_4, c_5

Distance matrix

$$D = \begin{pmatrix} c_1 & 0 & & & \\ c_2 & 9 & 0 & & \\ c_3 & 3 & 7 & 0 & \\ c_4 & 6 & 5 & 9 & 0 \\ c_5 & 11 & 10 & 2 & 8 & 0 \end{pmatrix}$$

$c_1 \ c_2 \ c_3 \ c_4 \ c_5$

$$d_{c_3 c_5} = \min d_{ij} = 2$$

\Rightarrow Merge c_3 & c_5 to form cluster (c_3, c_5) at merger level 2 - total 4 clusters

Row col^m deletion containing c_3 & c_5

$$\begin{pmatrix} 0 & & & & \\ 9 & 0 & & & \\ c_3 & 3 & 7 & 9 & - \\ c_5 & 6 & 5 & 0 & 1 \\ 11 & 10 & 2 & 8 & 0 \end{pmatrix}$$

$c_3 \quad c_5$

Computation of $d_{(c_3, c_5), c_i}$ for $i = 1, 2, 4$

$$\begin{aligned} d_{(c_3, c_5) c_1} &= \min(d_{c_1 c_3}, d_{c_1 c_5}) \\ &= \min(3, 11) = 3 \end{aligned}$$

$$d_{(c_3, c_5) c_2} = \min(d_{c_2 c_3}, d_{c_2 c_5}) = \min(7, 10) = 7$$

$$d_{(c_3, c_5) c_4} = \min(d_{c_3 c_4}, d_{c_4 c_5}) = \min(9, 8) = 8.$$

Updated distance matrix x

(C_3, C_5)	0				
C_1	3	0			
C_2	7	9	0		
C_4	8	6	5	0	

$$d_{(C_3, C_5)C_1} \text{ is min}$$

Merge (C_3, C_5) with C_1 to form cluster

(C_1, C_3, C_5) at merger level 3 - total 3 clusters

Row column deletion step

(C_3, C_5)	0	-	-	-	-
C_1	3	0	-	-	-
C_2	7	9	0		
C_4	8	6	5	0	

Calculate distance bet " (C_3, C_5, C_1) & C_2 and
 (C_3, C_5, C_1) & C_4

$$\begin{aligned} d_{(C_1, C_3, C_5)C_2} &= \min(d_{C_1C_2}, d_{(C_3, C_5)C_2}) \\ &= \min(9, 7) = 7 \end{aligned}$$

$$\begin{aligned} d_{(C_1, C_3, C_5)C_4} &= \min(d_{C_1C_4}, d_{(C_3, C_5)C_4}) \\ &= \min(6, 8) = 6 \end{aligned}$$

Update distance matrix

$$\begin{matrix} (c_1, c_3, c_5) & \left(\begin{array}{ccc} 0 & & \\ & & \\ & & \end{array} \right) \\ c_2 & \left(\begin{array}{cc} 7 & 0 \\ & \end{array} \right) \\ c_4 & \left(\begin{array}{ccc} 6 & 5 & 0 \\ & \boxed{5} & \\ & & \end{array} \right) \end{matrix}$$

d_{c_2, c_4} is min

Merge c_2 & c_4 to form cluster (c_2, c_4) at merger level 5 - total 2 clusters (c_1, c_3, c_5) & (c_2, c_4)

Row column deletion

$$\begin{matrix} (c_1, c_3, c_5) & \left(\begin{array}{cc|c} 0 & & \\ & & | \\ & & | \\ \hline c_2 & 7 & 0 \\ c_4 & 6 & 5 \end{array} \right) \\ c_2 & \left(\begin{array}{c} 7 \\ -0- \\ \hline \end{array} \right) \\ c_4 & \left(\begin{array}{c} 6 \\ -5- \\ \hline \end{array} \right) \end{matrix}$$

Distance matrix updation

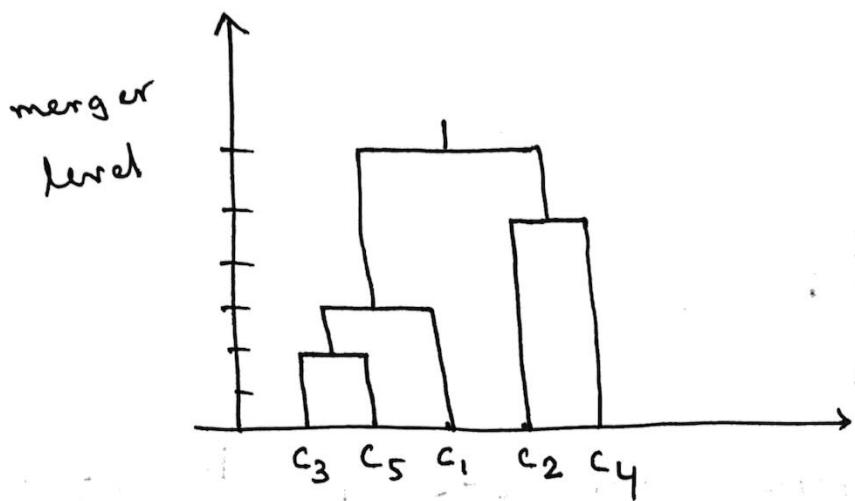
$$\begin{matrix} (c_1, c_3, c_5) & \left(\begin{array}{cc} 0 & \\ & \end{array} \right) \\ (c_2, c_4) & \left(\begin{array}{cc} ? & 0 \end{array} \right) \end{matrix}$$

$$\begin{aligned} d_{(c_1, c_3, c_5)(c_2, c_4)} &= \min(d_{(c_1, c_3, c_5)c_2}, d_{(c_1, c_3, c_5)c_4}) \\ &= \min(7, 6) = 6 \end{aligned}$$

$$\begin{matrix} (c_1, c_3, c_5) & \left(\begin{array}{cc} 0 & \\ & \end{array} \right) \\ (c_2, c_4) & \left(\begin{array}{cc} 6 & 0 \end{array} \right) \end{matrix}$$

\Rightarrow Merge (c_1, c_3, c_5) & (c_2, c_4) at merger level 6
- single cluster

Dendogram



Measures of dissimilarity

$$\mathbf{x} = (x_1 : \dots : x_n)$$

$\mathbf{x} \rightarrow D$: dissimilarity matrix
 $p \times p$

$$D = ((d_{rs})) \Rightarrow$$

$$(i) d_{rs} \geq 0 \quad \forall r, s$$

$$(ii) d_{rr} = 0 \quad \forall r$$

$$(iii) d_{rs} = d_{sr} \quad \forall r, s$$

Remark: (i) Symmetry cond" may not always be satisfied

$$D \rightarrow (D + D')/2$$

(ii) If in addition to the cond's (i)-(iii),
 d_{rs} also satisfies triangle inequality

$$d_{rs} \leq d_{rt} + d_{ts} \quad \forall r, s, t$$

then the dissimilarity measure is a metric and

We use the term "distance"

Some commonly used dissimilarity measures for numeric values

x_i, x_j : $p \times 1$ feature vector

$$x_i = \begin{pmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{pi} \end{pmatrix}$$

(a) Euclidean distance

$$d_e = \left(\sum_{i=1}^p (x_{1i} - x_{2i})^2 \right)^{1/2}$$

(b) Absolute value distance or City-block distance

$$d_{cb} = \sum_{i=1}^p |x_{1i} - x_{2i}|$$

(c) Chebyshov distance

$$d_{ch} = \max_i |x_{li} - x_{rj}|$$

(d) Pearson distance

$$d_p = \left(\sum_{i=1}^p (x_{li} - x_{rj})^2 / s_i^2 \right)^{1/2}$$

s_i^2 : sample variance for i^{th} variable

(e) Quadratic distance

$$d_q = \left((\underline{x}_i - \underline{x}_j)' Q (\underline{x}_i - \underline{x}_j) \right)^{1/2}$$

Q is p.d.

e.g. $Q = S^{-1}$: S = within group covariance matrix

case of feature vector with ordinal values

Realizable values are ordered set.

e.g. Academic grades: A, B, C, D, E, F

opinion ~~Opinion~~: strongly disagree, disagree, neutral, agree, strongly agree

One approach is to transform M (say) ordinal values

to

$$\frac{j - \frac{1}{2}}{M}; j = 1(1)M$$

in the prescribed order of their original values.

Treat the transformed values as quantitative numeric values.

Case of nominal (categorical, unordered) / ordinal values

Transform to binary variables set.

Nominal variable with K states

→ represent as K binary variables

nominal variable at m^{th} state

$$(0, \dots, \underset{m^{\text{th}}}{1}, \dots, 0)$$

SLy for ordinal

Similarity measures for binary variables

Let \underline{x} & \underline{y} be 2 vectors of binary variables.

Define

a : # of occurrences of $x_i = 1$ and $y_i = 1$

b : # of occurrences of $x_i = 0$ & $y_i = 1$

c : $x_i = 1$ & $y_i = 0$

d : $x_i = 0$ & $y_i = 0$

$a + b + c + d = p$, # of variables.

Some similarity measures:

(a) Simple matching coefficient:

$$\delta_{sm} = \frac{a+d}{a+b+c+d}$$

(b) Russell & Rao

$$\delta_{RR} = \frac{a}{a+b+c+d}$$

(c) Jaccard

$$\delta_J = \frac{a}{a+b+c}$$

(d) Czekanowski:

$$\delta_{CZ} = \frac{2a}{2a+b+c}$$

Note:

dissimilarity measure

$$\text{e.g. } d_{xy} = 1 - \delta_{xy} \text{ for } S_m$$

Mixed type: dimensions have both numerical / categorical

Transform all to binary.

Example

	Height*	Weight*	Eye color	Hair color	Handedness	Gender
Ind 1	68	140	green	blond	R	F
Ind 2	73	185	brown	brown	R	M
Ind 3	67	165	blue	blond	R	M
Ind 4	64	120	brown	brown	R	F
Ind 5	76	210	brown	brown	L	M

* : in inch

* : in lb

Define, $x_{1i} = \begin{cases} 1, & ht \geq 72 \text{ in} \\ 0, & \text{o/w.} \end{cases}$ $x_{2i} = \begin{cases} 1, & wt \geq 150 \text{ lb} \\ 0, & \text{o/w} \end{cases}$

$$x_{3i} = \begin{cases} 1, & \text{brown eyes} \\ 0, & \text{o/w} \end{cases} \quad x_{4i} = \begin{cases} 1, & \text{blond hair} \\ 0, & \text{o/w} \end{cases}$$

$$x_{5i} = \begin{cases} 1, & R \\ 0, & L \end{cases} \quad x_{6i} = \begin{cases} 1, & \text{female} \\ 0, & \text{male} \end{cases}$$

	x_1	x_2	x_3	x_4	x_5	x_6
Ind 1	0	0	0	1	1	1
Ind 2	1	1	1	0	1	0

 \rightarrow

	Ind 1		
Ind 2	1	0	
	1 (a)	2 (b)	
0	3 (c)	0 (d)	

$$s_{dm}^{(1,2)} = \frac{1}{6}$$

Combine $\neq (i, j)$ $i \neq j$ pairs and get similarity matrix

$$S = \begin{pmatrix} 1 & 1 & & & & \\ 2 & \frac{1}{6} & 1 & & & \\ 3 & - & - & 1 & & \\ 4 & - & - & - & 1 & \\ 5 & - & - & - & - & 1 \\ 6 & - & - & - & - & - \end{pmatrix}$$

Distance betⁿ groups of objects

(A) Methods based on feature vectors

Already discussed (single linkage, complete linkage, average linkage, centroid distance, median distance)

(B) Methods based on probabilistic distance.

Distance measures use complete information about the structure of groups (classes) provided by the class conditional densities. The distance measure, say \mathcal{J} , satisfies

$$(i) \quad \mathcal{J} = 0 \quad \text{if} \quad p(x|\pi_1) = p(x|\pi_2)$$

$$(ii) \quad \mathcal{J} \geq 0$$

(iii) \mathcal{J} attains the maximum value when the classes are disjoint

Commonly used probabilistic distances

(I) Bhattacharya distance

$$\mathcal{J}_B = -\log \left(\int \{p(x|\pi_1) p(x|\pi_2)\}^{\frac{1}{2}} dx \right)$$

(II) Chernoff distance

$$\mathcal{J}_C = -\log \left(\int \{p(x|\pi_1) \hat{p}(x|\pi_2)\}^{1-\delta} dx \right)$$

$$\delta \in [0, 1]$$

(III) Divergence distance

$$\mathcal{J}_D = \int (p(x|\pi_1) - p(x|\pi_2)) \log \left(\frac{p(x|\pi_1)}{p(x|\pi_2)} \right) dx$$

(iv) Patrick - Fischer distance

$$J_{PF} = \left(\int \left\{ p_1 p(\underline{x} | \pi_1) - p_2 p(\underline{x} | \pi_2) \right\}^2 d\underline{x} \right)^{1/2}$$

p_1 & p_2 are prior probabilities of the classes.

Remark: Probabilistic distances are computed under specific distributional setups.

Example: Gtr I objects - $N_p(\underline{\mu}_1, \Sigma_1)$

$$\Sigma_i > 0$$

Gtr II objects - $N_p(\underline{\mu}_2, \Sigma_2)$

Divergence distance calculations

$$J_D = \int \left((2\pi)^{-p/2} |\Sigma_1|^{-1/2} \exp\left(-\frac{1}{2}(\underline{x} - \underline{\mu}_1)' \Sigma_1^{-1} (\underline{x} - \underline{\mu}_1)\right) \right. \\ \left. - (2\pi)^{-p/2} |\Sigma_2|^{-1/2} \exp\left(-\frac{1}{2}(\underline{x} - \underline{\mu}_2)' \Sigma_2^{-1} (\underline{x} - \underline{\mu}_2)\right) \right) \\ \left(\log \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} + \log \left(\exp\left(-\frac{1}{2}((\underline{x} - \underline{\mu}_1)' \Sigma_1^{-1} (\underline{x} - \underline{\mu}_1) - (\underline{x} - \underline{\mu}_2)' \Sigma_2^{-1} (\underline{x} - \underline{\mu}_2))\right) \right) \right) d\underline{x}$$

$$= \log \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} \int (p(\underline{x} | \pi_1) - p(\underline{x} | \pi_2)) d\underline{x}$$

$$+ \int \left(-\frac{1}{2} (\underline{x} - \underline{\mu}_1)' \Sigma_1^{-1} (\underline{x} - \underline{\mu}_1) \right) \left(p(\underline{x} | \pi_1) - p(\underline{x} | \pi_2) \right) d\underline{x}$$

$$+ \int \left(\frac{1}{2} (\underline{x} - \underline{\mu}_2)' \Sigma_2^{-1} (\underline{x} - \underline{\mu}_2) \right) \left(p(\underline{x} | \pi_1) - p(\underline{x} | \pi_2) \right) d\underline{x}$$

$$= (1 - 1) \log \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}}$$

$$+ \left(-\frac{1}{2}\right) \int (\underline{x}' \Sigma_1^{-1} \underline{x} + \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 - 2 \underline{\mu}_1' \Sigma_1^{-1} \underline{x}) \\ (\Pr(\underline{x} | \pi_1) - \Pr(\underline{x} | \pi_2)) d\underline{x}$$

$$+ \left(\frac{1}{2}\right) \int (\underline{x}' \Sigma_2^{-1} \underline{x} + \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 - 2 \underline{\mu}_2' \Sigma_2^{-1} \underline{x}) \\ (\Pr(\underline{x} | \pi_1) - \Pr(\underline{x} | \pi_2)) d\underline{x}$$

$$= -\frac{1}{2} \left(\left\{ E_{\pi_1} (\underline{x}' \Sigma_1^{-1} \underline{x}) + \cancel{\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1} - 2 \underline{\mu}_1' \Sigma_1^{-1} E_{\pi_1} (\underline{x}) \right\} \right. \\ \left. - \left\{ E_{\pi_2} (\underline{x}' \Sigma_1^{-1} \underline{x}) + \cancel{\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1} - 2 \underline{\mu}_1' \Sigma_1^{-1} E_{\pi_2} (\underline{x}) \right\} \right)$$

$$+ \frac{1}{2} \left(\left\{ E_{\pi_1} (\underline{x}' \Sigma_2^{-1} \underline{x}) + \cancel{\underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2} - 2 \underline{\mu}_2' \Sigma_2^{-1} E_{\pi_1} (\underline{x}) \right\} \right.$$

$$\left. - \left\{ E_{\pi_2} (\underline{x}' \Sigma_2^{-1} \underline{x}) + \cancel{\underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2} - 2 \underline{\mu}_2' \Sigma_2^{-1} E_{\pi_2} (\underline{x}) \right\} \right)$$

— *

Note that

$$(i) E_{\pi_i} (\underline{x}) = \underline{\mu}_i ; i = 1, 2$$

and we need to calculate terms of type

$$E_{\pi_i} (\underline{x}' \Sigma_j^{-1} \underline{x}) ; i, j = 1, 2 \\ = ?$$

If $\underline{x} \sim N_p(\underline{\mu}, \Sigma)$, then

$$\begin{aligned}
 E(\underline{x}' \Sigma^{-1} \underline{x}) &= E(\text{tr}(\underline{x}' \Sigma^{-1} \underline{x})) \\
 &= E(\text{tr}(\Sigma^{-1} \underline{x} \underline{x}')) \\
 &= \text{tr}(E(\Sigma^{-1} \underline{x} \underline{x}')) \\
 &= \text{tr}(\Sigma^{-1} E(\underline{x} \underline{x}')) \\
 &= \text{tr}(\Sigma^{-1} (\Sigma + \underline{\mu} \underline{\mu}')) \\
 &= \text{tr} \Sigma^{-1} \Sigma + \text{tr}(\Sigma^{-1} \underline{\mu} \underline{\mu}') \\
 &= \text{tr} I_p + \text{tr}(\underline{\mu}' \Sigma^{-1} \underline{\mu}) \\
 &= \text{tr} I_p + \underline{\mu}' \Sigma^{-1} \underline{\mu} \quad - (**)
 \end{aligned}$$

In general,

$$\text{From } (*), E(\underline{x}' A \underline{x}) = \text{tr} A \Sigma + \underline{\mu}' A \underline{\mu}$$

$$\begin{aligned}
 J_D &= (\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 - \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_2 \\
 &\quad - \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_1 + \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2) \\
 &\quad - \frac{1}{2} \left(E_{\pi_1}(\underline{x}' \Sigma_1^{-1} \underline{x}) - E_{\pi_2}(\underline{x}' \Sigma_1^{-1} \underline{x}) \right. \\
 &\quad \left. - E_{\pi_1}(\underline{x}' \Sigma_2^{-1} \underline{x}) + E_{\pi_2}(\underline{x}' \Sigma_2^{-1} \underline{x}) \right)
 \end{aligned}$$

$$\begin{aligned}
&= \left(\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 + \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 - \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_2 - \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_1 \right) \\
&\quad - \frac{1}{2} \left(\left(\text{tr} \Sigma_1^{-1} \Sigma_1 + \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 \right) - \left(\text{tr} \Sigma_1^{-1} \Sigma_2 + \underline{\mu}_2' \Sigma_1^{-1} \underline{\mu}_2 \right) \right. \\
&\quad \left. - \left(\text{tr} \Sigma_2^{-1} \Sigma_1 + \underline{\mu}_1' \Sigma_2^{-1} \underline{\mu}_1 \right) + \left(\text{tr} \Sigma_2^{-1} \Sigma_2 + \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 \right) \right) \\
&= \left(\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 + \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 - \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_2 - \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_1 \right. \\
&\quad \left. - \frac{1}{2} \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 + \frac{1}{2} \underline{\mu}_2' \Sigma_1^{-1} \underline{\mu}_2 + \frac{1}{2} \underline{\mu}_1' \Sigma_2^{-1} \underline{\mu}_1 - \frac{1}{2} \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 \right) \\
&\quad + \left(-\frac{1}{2} \text{tr} I_p + \frac{1}{2} \text{tr} \Sigma_1^{-1} \Sigma_2 + \frac{1}{2} \text{tr} \Sigma_2^{-1} \Sigma_1 - \frac{1}{2} \text{tr} I_p \right) \\
&= \left(\frac{1}{2} \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 + \frac{1}{2} \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 + \frac{1}{2} \underline{\mu}_2' \Sigma_1^{-1} \underline{\mu}_2 + \frac{1}{2} \underline{\mu}_1' \Sigma_2^{-1} \underline{\mu}_1 \right. \\
&\quad \left. - \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 - \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_2 \right) \\
&\quad + \frac{1}{2} \left(\text{tr} \Sigma_1^{-1} \Sigma_2 + \text{tr} \Sigma_2^{-1} \Sigma_1 - 2 \text{tr} I_p \right) \\
&= \frac{1}{2} \left(\underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_1 + \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_2 + \underline{\mu}_2' \Sigma_1^{-1} \underline{\mu}_2 + \underline{\mu}_1' \Sigma_2^{-1} \underline{\mu}_1 \right. \\
&\quad \left. - 2 \underline{\mu}_1' \Sigma_1^{-1} \underline{\mu}_2 - 2 \underline{\mu}_2' \Sigma_2^{-1} \underline{\mu}_1 \right) \\
&\quad + \frac{1}{2} \left(\text{tr} \Sigma_1^{-1} \Sigma_2 + \text{tr} \Sigma_2^{-1} \Sigma_1 - 2 \text{tr} I_p \right) \\
&= \frac{1}{2} (\underline{\mu}_2 - \underline{\mu}_1)' (\Sigma_1^{-1} + \Sigma_2^{-1}) (\underline{\mu}_2 - \underline{\mu}_1) \\
&\quad + \frac{1}{2} \text{tr} (\Sigma_1^{-1} \Sigma_2 + \Sigma_2^{-1} \Sigma_1 - 2 I_p).
\end{aligned}$$

Note: If $\Sigma_1 = \Sigma_2$, then

$$J_D = (\underline{\mu}_2 - \underline{\mu}_1)' \Sigma^{-1} (\underline{\mu}_2 - \underline{\mu}_1)$$

→ square of Mahalanobis distance

Note: For $N_p(\underline{\mu}_1, \Sigma_1)$ & $N_p(\underline{\mu}_2, \Sigma_2)$ setup.

Chernoff distance

$$J_C = \frac{1}{2} \lambda(1-\lambda) (\underline{\mu}_2 - \underline{\mu}_1)' \Sigma_{\lambda}^{-1} (\underline{\mu}_2 - \underline{\mu}_1) + \frac{1}{2} \log \left(\frac{|\Sigma_{\lambda}|}{|\Sigma_1|^{1-\lambda} |\Sigma_2|^{\lambda}} \right)$$

$$\Sigma_{\lambda} = (1-\lambda) \Sigma_1 + \lambda \Sigma_2 ; \lambda \in [0, 1]$$

Note: If $\Sigma_1 = \Sigma_2$ i.e. $N_p(\underline{\mu}_1, \Sigma)$ & $N_p(\underline{\mu}_2, \Sigma)$.

$$J_D = 8 J_B = (\underline{\mu}_2 - \underline{\mu}_1)' \Sigma^{-1} (\underline{\mu}_2 - \underline{\mu}_1)$$

Note: For $N_p(\underline{\mu}_1, \Sigma_1)$ & $N_p(\underline{\mu}_2, \Sigma_2)$ setup

$$J_{PF} = (2\pi)^{-p/2} \left(|\Sigma_1|^{-1/2} + |\Sigma_2|^{-1/2} - 2 |\Sigma_1 + \Sigma_2|^{-1/2} \right)$$

Patrick-Fischer

$$\exp \left(-\frac{1}{2} (\underline{\mu}_2 - \underline{\mu}_1)' (\Sigma_1 + \Sigma_2)^{-1} (\underline{\mu}_2 - \underline{\mu}_1)' \right)$$

Non-hierarchical clustering methods

Features

- no distance/dissimilarity matrix calculations
- no hierarchy of clusters
- starts from either
 - (i) a random initial partition of items into groups
 - or (ii) an initial set of randomly selected seed points, which will form the nuclei of clusters.

K-means clustering or iterative relocation method

K-means method is an iterative algorithm that assigns each item to the cluster having the nearest centroid (mean).

Steps for K-means algorithm

- (I) Partition the items into K initial clusters.
- (II) Reassign items to the cluster whose centroid is nearest (in Euclidean sense mostly). Recalculate the centroids for the cluster receiving the new item and for the cluster losing that item.
- (III) Repeat step (II) till no more reassignment can take place.

Note: Rather than starting with a partition of items into K initial groups in step (I), we can also specify

K initial seed points (centroids) and then proceed to step (ii) after a walk through the data.

Remark: Let's try to understand the logic behind K-means method !!

Let there be N objects to be put into K clusters.

Suppose $c(i)=k$ denote that object i in cluster k

$$k=1(1)K; i=1(1)N$$

A natural "loss function" (criterion) for clustering :

$$W(c) = \frac{1}{2} \sum_{k=1}^{K} \sum_{c(i)=k} \sum_{c(i')=k} d_{ii'} - (1)$$

$$\text{where, } d_{ii'} = d(x_i, x_{i'})$$

$W(c)$ is within cluster point scatter, a measure of compactness of clusters. It characterizes the extent to which observations assigned, under a partition, to same clusters tend to be close to one another.

Further, consider the total point scatter

$$T = \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N d_{ii'}$$

$$\text{i.e. } T = \frac{1}{2} \sum_{k=1}^{K} \sum_{c(i)=k} \left(\sum_{c(i')=k} d_{ii'} + \sum_{c(i') \neq k} d_{ii'} \right)$$

$$\text{i.e. } T = \frac{1}{2} \sum_{k=1}^K \sum_{c(i)=k} \sum_{c(i')=k} d_{ii'} + \frac{1}{2} \sum_{k=1}^K \sum_{c(i)=k} \sum_{c(i') \neq k} d_{ii'}$$

$$\text{i.e. } T = W(c) + B(c)$$

$$B(c) = \frac{1}{2} \sum_{k=1}^K \sum_{c(i)=k} \sum_{c(i') \neq k} d_{ii'} \quad \text{is between-}$$

cluster point scatter. A measure of separation of clusters; large when different clusters are far apart.

Note:

$W(c)$ & $B(c)$ depends on cluster assignment
 T is indep of cluster assignment

Note: It is natural to maximize $B(c)$ or equivalently minimize $W(c)$ over all possible cluster assignments

Note: The combinatorial optimization is not feasible even for moderate N & K .

The number of distinct cases (assignments) is

$$S(N, K) = \frac{1}{K!} \sum_{k=1}^K (-1)^{K-k} \binom{K}{k} k^N$$

$$\text{e.g } S(10, 4) = 34,105$$

$$S(19, 4) \approx 10^{10}$$

Feasible strategies are based on iterative greedy descent - K-means algorithm is one such algorithm.

K-means algorithm

Iterative descent clustering algorithm

Let

$$d(\underline{x}_i, \underline{x}_{i'}) = d_{ii'} = \sum_{j=1}^p (x_{ij} - x_{i'j})^2 = \|\underline{x}_i - \underline{x}_{i'}\|^2$$

Within-cluster point scatter

$$\begin{aligned} W(c) &= \frac{1}{2} \sum_{k=1}^K \sum_{c(i)=k} \sum_{c(i')=k} \|\underline{x}_i - \underline{x}_{i'}\|^2 \\ &= \sum_{k=1}^K N_k \sum_{c(i)=k} \|\underline{x}_i - \bar{\underline{x}}_k\|^2 \end{aligned}$$

where $N_k = \sum_{i=1}^N I(c(i)=k)$

$\bar{\underline{x}}_k$: mean of of k^{th} cluster

Objective is thus to get $c^* \ni$

$$c^* = \operatorname{argmin}_c \sum_{k=1}^K N_k \sum_{c(i)=k} \|\underline{x}_i - \bar{\underline{x}}_k\|^2$$

and we can obtain c^* by enlarged optimization problem

$$\min_{c, \underline{m}_1, \dots, \underline{m}_K} \sum_{k=1}^K N_k \sum_{c(i)=k} \|\underline{x}_i - \underline{m}_k\|^2 \quad - (2)$$

as for any set of S observations

$$\bar{\underline{x}}_S = \operatorname{argmin}_{\underline{m}} \sum_{i=1}^S \|\underline{x}_i - \underline{m}\|^2$$

K-mean algorithm is an alternating optimization procedure for minimization of (2)

K-means clustering algorithm

Step I. For a given cluster assignment C , the total within cluster variability:

$$\sum_{k=1}^K N_k \sum_{\{c(i)=k\}} \|x_i - \bar{x}_k\|^2$$

is minimized w.r.t. $\bar{x}_1, \dots, \bar{x}_K$ by taking means of the currently assigned clusters

Step II: Given a current set of cluster means,

$$\sum_{k=1}^K N_k \sum_{\{c(i)=k\}} \|x_i - \bar{x}_k\|^2 \quad \text{i.e. } \bar{x}_k = \bar{x}_k$$

is minimized by assigning each obsn to the closest cluster mean

$$\text{i.e. } C(i) = \underset{1 \leq k \leq K}{\operatorname{argmin}} \|x_i - \bar{x}_k\|^2$$

Step I & II are iterated until the assignments do not change.

Remark: The result may represent a suboptimal local minimum.

One should start the algorithm with many initial random choices and choose the solⁿ having minimum objective function value.

Example : K-means

Cases	Variables	
	x_1	x_2
A	5	3
B	-1	1
C	1	-2
D	-3	-2

$K=2$

Step I : Arbitrary partition (A, B) (C, D)

Cluster centroids

Cluster	\bar{x}_1	\bar{x}_2
(A, B)	2	2
(C, D)	-1	-2

Step II : Compute distances from cluster centroids and
reassign each item to the nearest group

(If an item is moved from the initial configuration,
the cluster centroid must be updated before we
can proceed further)

$$\underline{A} \quad d^2(A, (A, B)) = 3^2 + 1^2 = 10$$

$$d^2(A, (C, D)) = 6^2 + 5^2 = 61$$

$\Rightarrow A$ is closer to (A, B) than (C, D)

\Rightarrow No reassignment required.

$$\underline{B} \quad d^2(B, (A, B)) = 3^2 + 1^2 = 10$$

$$d^2(B, (C, D)) = 0^2 + 3^2 = 9$$

\Rightarrow Reassign B to (c, D) to form (B, c, D)

(cluster centroid updation)

Cluster	Coordinates of centroid	
	\bar{x}_1	\bar{x}_2
A	5	3
(B, c, D)	-1	-1

$$d^2(A, A) = 0 ; d^2(A, (B, c, D)) = 52$$

$$d^2(B, A) = 40 ; d^2(B, (B, c, D)) = 4$$

$$d^2(c, A) = 41, d^2(c, (B, c, D)) = 5$$

$$d^2(D, A) = 89, d^2(D, (B, c, D)) = 5$$

\Rightarrow no reassignment is necessary and relocation iteration stops.

(A), (B, c, D) is the final clustering.

Comparing different cluster partitions

Objective is to have a criterion function for comparing different cluster partitions of the data and hence to choose the best cluster partition.

Let the N data points be given by $\tilde{x}_1, \dots, \tilde{x}_N$.

The sample variance covariance matrix $\hat{\Sigma}$ is given by

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\tilde{x}_i - \tilde{m})(\tilde{x}_i - \tilde{m})'$$

$$\text{where } \tilde{m} = \frac{1}{N} \sum_{i=1}^N \tilde{x}_i = \bar{\tilde{x}}$$

Let there be K clusters and define

$$z_{ji} = \begin{cases} 1, & \text{if } \tilde{x}_i \in \text{cluster } j \\ 0, & \text{otherwise} \end{cases}$$

$$\tilde{m}_j = \frac{1}{N_j} \sum_{i=1}^{N_j} z_{ji} \tilde{x}_i (= \bar{\tilde{x}}_j) \quad \text{mean of cluster } j$$

$$n_j = \sum_{i=1}^{N_j} z_{ji} \quad \# \text{ of } \tilde{x}_i \text{'s in cluster } j$$

The within-cluster sum of squares and cross product (SSCP) scatter matrix is

$$S_W = \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^{N_j} z_{ji} (\tilde{x}_i - \tilde{m}_j)(\tilde{x}_i - \tilde{m}_j)'$$

(The pooled within cluster scatter matrix over K clusters)

Note that

$$\begin{aligned}
 \sum &= \frac{1}{N} \sum_{i=1}^N (\tilde{x}_i - \tilde{m})(\tilde{x}_i - \tilde{m})' \\
 &= \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\tilde{x}_i - \tilde{m})(\tilde{x}_i - \tilde{m})' \\
 &= \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\overline{\tilde{x}_i - \tilde{m}_j} + \overline{\tilde{m}_j - \tilde{m}})(\tilde{x}_i - \tilde{m}_j + \tilde{m}_j - \tilde{m})' \\
 &= \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} \left((\tilde{x}_i - \tilde{m}_j)(\tilde{x}_i - \tilde{m}_j)' \right. \\
 &\quad \left. + (\tilde{m}_j - \tilde{m})(\tilde{m}_j - \tilde{m})' \right. \\
 &\quad \left. + \cancel{(\tilde{x}_i - \tilde{m}_j)(\tilde{m}_j - \tilde{m})'} \right. \\
 &\quad \left. + \cancel{(\tilde{m}_j - \tilde{m})(\tilde{x}_i - \tilde{m}_j)'} \right)
 \end{aligned}$$

i.e. $\sum = \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\tilde{x}_i - \tilde{m}_j)(\tilde{x}_i - \tilde{m}_j)'$

$$+ \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\tilde{m}_j - \tilde{m})(\tilde{m}_j - \tilde{m})'$$

(e.g. $\frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\tilde{m}_j - \tilde{m})(\tilde{x}_i - \tilde{m}_j)'$)

$$= \frac{1}{N} \sum_{j=1}^K (\tilde{m}_j - \tilde{m}) \sum_{i=1}^N z_{ji} (\tilde{x}_i - \tilde{m}_j)'$$

$$= \frac{1}{N} \sum_{j=1}^K (\tilde{m}_j - \tilde{m}) (n_j \tilde{m}_j - n_j \tilde{m}_j)' = 0$$

So $\frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\tilde{x}_i - \tilde{m}_j)(\tilde{m}_j - \tilde{m})' = 0$

i.e.

$$\sum = S_W + \frac{1}{N} \sum_{j=1}^K n_j (\underline{m}_j - \underline{m})(\underline{m}_j - \underline{m})'$$

$$= S_W + S_B$$

S_B : between-cluster sum of squares and cross product
 (SSCP) \wedge scatter matrix
 (indicates the scatter of the cluster means
 about total grand mean)

Popular optimum clustering criteria are based on
 univariate (scalar) functions of above matrices, S_W, S_B, \sum

e.g. (a) Minimization of $\text{tr}(S_W)$

$$\begin{aligned}\text{tr}(S_W) &= \text{tr} \left(\frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} (\underline{x}_i - \underline{m}_j)(\underline{x}_i - \underline{m}_j)' \right) \\ &= \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} \text{tr}((\underline{x}_i - \underline{m}_j)(\underline{x}_i - \underline{m}_j)') \\ &= \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} \text{tr}((\underline{x}_i - \underline{m}_j)'(\underline{x}_i - \underline{m}_j)) \\ &= \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N z_{ji} \|\underline{x}_i - \underline{m}_j\|^2\end{aligned}$$

i.e. $\text{tr}(S_W) = \frac{1}{N} \sum_{j=1}^K S_j$

$$S_j = \sum_{i=1}^N z_{ji} \|\underline{x}_i - \underline{m}_j\|^2$$

The within-n-cluster sum of squares for j^{th} cluster

Thus minimization of $\text{tr}(S_W)$ is same as
minimization of total within cluster sum of squares
about the K centroids (equiv to max $\text{tr} S_B$)

$$(b) \text{ Min } \frac{\|S_W\|}{\|\Sigma\|}$$

↳ Solution of eqn 1 has been mentioned above in handwritten notes

$$(c) \text{ Min } \text{tr}(\Sigma^{-1} S_W)$$

↳ Solution of eqn 2 has been mentioned above in handwritten notes

$$(d) \text{ Max } \text{tr}(S_W^{-1} S_B)$$

↳ Solution of eqn 3 has been mentioned above in handwritten notes

↳ Solution of eqn 4 has been mentioned above in handwritten notes

↳ Solution of eqn 5 has been mentioned above in handwritten notes

↳ Solution of eqn 6 has been mentioned above in handwritten notes

↳ Solution of eqn 7 has been mentioned above in handwritten notes

↳ Solution of eqn 8 has been mentioned above in handwritten notes

↳ Solution of eqn 9 has been mentioned above in handwritten notes

↳ Solution of eqn 10 has been mentioned above in handwritten notes

↳ Solution of eqn 11 has been mentioned above in handwritten notes

↳ Solution of eqn 12 has been mentioned above in handwritten notes