

Cluster Analysis

(1) Similar Coefficient:

① Transf of Data:

If we have data $X = (x_{ij})_{n \times m}$, $(x_{ij})_{j=1}^m$ is from π_i . Then we obtain:

$$\text{mean: } \bar{x}_j = \sum_i x_{ij} / n, \quad \text{SD: } s_j = \sqrt{\frac{1}{n-1} \sum_i (x_{ij} - \bar{x}_j)^2}$$

$$\text{Range: } R_j = \max_i x_{ij} - \min_i x_{ij}.$$

i) Centralization:

$$x_{ij}^* = x_{ij} - \bar{x}_j, \quad \text{i.e. } X^* = X - \frac{1}{n} J_n^n X.$$

$$S^* = \sum_{i=1}^n x_{ti}^* x_{ti}^* / (n-1) = S \quad (\text{invariant})$$

ii) Standardization:

$$x_{ij}^* = (x_{ij} - \bar{x}_j) / s_j, \quad \text{indep't with unit.}$$

iii) Range Standardization:

$$x_{ij}^* = (x_{ij} - \bar{x}_j) / R_j, \quad \text{indep't with unit.}$$

$$\text{Rank Normalization: } x_{ij}^* = (x_{ij} - \min_i x_{ij}) / R_j$$

iv) logarithm:

$$x_{ij}^* = \log x_{ij}, \quad \text{where } x_{ij} > 0. \quad \text{For linearizing.}$$

Def: i) Minkowski Dist: $d_{ij}(p) = \left(\sum_i |x_{it} - x_{jt}|^p \right)^{\frac{1}{p}}$

ii) Lance Dist: $d_{ij}(L) = \frac{1}{m} \sum_i |x_{it} - x_{jt}| / (x_{it} + x_{jt})$

iii) Bias Dist: $d_{ij} = \frac{1}{m} \left(\sum_k \sum_l (x_{ik} - x_{jk})(x_{il} - x_{jl}) r_{kl} \right)^{\frac{1}{2}}$

Rmk: i), ii) are often used in indept r.v.'s but
iii) isn't. ii) isn't sensitive with large eigenvalue.

Def: 在数量化理论, 将定性变量称为项目, 而
项目中不同的值称为类目.

e.g. 若有 n 个样品 X_i ($1 \leq i \leq n$), m 个项目.

k^{th} 项目有 r_k 个类目, 则定义 X_i 的取值:

$$X_i = (\delta_i(k, 1), \dots, \delta_i(k, r_k)), \quad 1 \leq k \leq m.$$

$$\delta_i(k, l) = \begin{cases} 1, & i^{th} \text{ sample 属于 } k^{th} \text{ 项目 } l^{th} \text{ 类目.} \\ 0, & \text{otherwise.} \end{cases}$$

Def: X_i couples with X_j on k^{th} 项目 l^{th} 类目 if

$$\delta_i(k, l) = \delta_j(k, l), \text{ 若都为 } 1 \text{ 则称为 } 1-1 \text{ 配对}$$

对. 若都为 0, 则称为 0-0 配对.

Def: In qualitative case, Distance between X_i
and X_j is defined by:

i) m_1 为 X_i, X_j 的 1-1 配对数; m_0 为 0-0 配

对数, m_2 为不配对数. ($m_1 + m_2 + m_0 = \sum_{k=1}^m r_k$)

$$d_{ij} = m_2 / (m_1 + m_2)$$

ii) $d_{ij}^2 = \sum_k \sum_l (\delta_i(k, l) - \delta_j(k, l))^2$, i.e. 不

配对的总数.

Rmk: 存在情况: sample X_i 属于 k^{th} 项目
不同的项目, 若不能同时兼取时.

Define distance in i : m_i^* 为不配
对项目总数. $d_{ij} = m_i^* / m$.

Def: Similar coefficient between X_i, X_j is C_{ij} .

satisfies i) $|C_{ij}| \leq 1, \forall i, j$ ii) $C_{ij} = \pm 1 \Leftrightarrow X_i = aX_j, a \neq 0$.

iii) $C_{ij} = C_{ji}$ iv) $|C_{ij}| \uparrow, X_i, X_j$ are more similar.

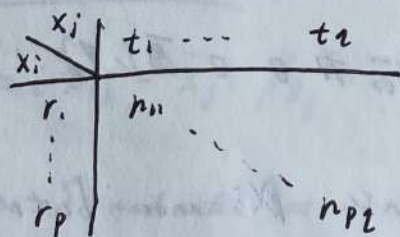
eg. i) $C_{ij}(1) = \frac{X_i^T X_j}{\sqrt{(X_i^T X_i)(X_j^T X_j)}} \quad (\cos \theta)$

ii) $C_{ij}(2) = r_{ij} = \frac{\sum (X_{ti} - \bar{X}_i)(X_{tj} - \bar{X}_j)}{\sqrt{\sum (X_{ti} - \bar{X}_i)^2 \sum (X_{tj} - \bar{X}_j)^2}}$

Rmk: Define distance by C_{ij} :

$d_{ij} = 1 - |C_{ij}|$ or $\hat{d}_{ij} = 1 - C_{ij}$.

iii) In qualitative case: If X_i 有 p 个项目, 取值



为 (r_1, \dots, r_p) , X_j 有 2 个项目
取值为 (t_1, \dots, t_2) .

n_{kl} = number of (X_i, X_j)
 $= (r_k, t_l)$

$\chi_{ij}^2 = n \left(\sum_{k,l} \left(\frac{n_{kl}^2}{n_{k \cdot} n_{\cdot l}} \right) - 1 \right)$

Def: $C_{ij} = \sqrt{\chi_{ij}^2 / (\chi_{ij}^2 + n)}$

(2) Hierarchical Method:

If we have samples $X_i = (X_{i1}, \dots, X_{im})$, $1 \leq i \leq n$.

The ideal is seeing n samples as n classes initially. Then define a distance. Classify two sample as a class which has minimal distance. Repeat the procedure...

Process: i) Data Transform ii) Calculate the matrix of distance $(d_{ij})_{n \times n} \stackrel{\Delta}{=} D_0$.

iii) If $(k, l) = \arg \min_{i,j} d_{ij}$. Then merge X_k, X_l as a class $\{X_k, X_l\}$.

iv) Calculate the distance w.r.t remaining classes to obtain $(d'_{ij})_{m \times m} \stackrel{\Delta}{=} D_1$.

v) Repeat the procedure until remain one class.

e.g. $X = (X_1, \dots, X_5) = (1, 2, 4, 5, 6, 8)$. $d_{ij} = |X_i - X_j|$.

i) $D^0 = \begin{pmatrix} 0 & \boxed{1} & 3.5 & 5 & 7 \\ 0 & 0 & 2.5 & 4 & 6 \\ & 0 & 0 & 1.5 & 3.5 \\ & & & 0 & 2 \\ & & & & 0 \end{pmatrix}$ $G_i = \{X_i\}$, $1 \leq i \leq 5$
 $k=5$

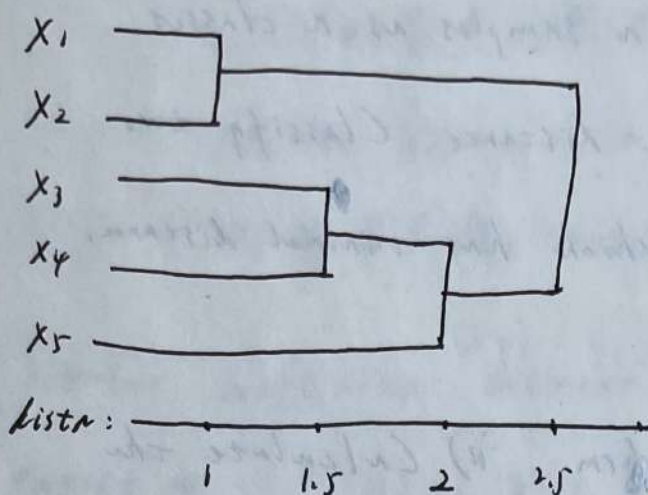
ii) d_{12} is min \Rightarrow Merge X_1, X_2 as $CL4 = \{X_1, X_2\}$

iii) $D^1 = \begin{array}{c|cccc} & X_3 & X_4 & X_5 & CL4 \\ \hline X_3 & & & & \\ X_4 & & & & \\ X_5 & & & & \\ CL4 & & & & \end{array}$ $d_{ij} = |X_i - X_j|$
 $d'_{ij} = \min_{X_j \in CL4} d(X_i, X_j)$

d'_{12} is min $\Rightarrow CL3 = \{X_3, X_4\}$

iv) $D^2 = \begin{array}{c|ccc} & X_5 & CL4 & CL3 \\ \hline X_5 & & & \\ CL4 & & & \\ CL3 & & & \end{array}$ $\Rightarrow CL2 = \{X_5, X_3, X_4\}$

v) Finally. $CL1 = \{CL2, CL4\}$.



Hierarchical Graph

iv) We can determine the number of groups by requiring the distance of group $< d_0$.
Ex. If require $d < d_0 = 2$. Then classify:

$$\{X_1, X_2\}, \{X_3, X_4\}, \{X_5\}, \quad k=3.$$

$$d < 3 : \{X_1, X_2\}, \{X_3, X_4, X_5\}, \quad k=2.$$

(i) Distance of Groups:

$$i) D_{12}^L = \mathcal{L}(G_1, G_2) =: \min \{d_{ij} \mid X_i \in G_1, X_j \in G_2\}$$

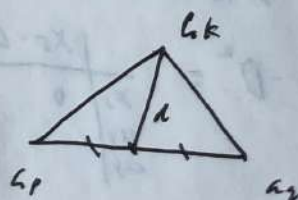
$$\text{or } D_{12}^S = \max \{d_{ij} \mid X_i \in G_1, X_j \in G_2\}$$

Rule: If merge G_1, G_2 to G_r . Then $D_{rk} = \min \{D_{pk}, D_{2k}\}$ or $= \max \{D_{pk}, D_{2k}\}$.

$$ii) \text{ By Formula: } D_{rk}^2 = \frac{1}{2} (D_{1k}^2 + D_{2k}^2) + \beta D_{12}^2, \quad -\frac{1}{4} \leq \beta \leq 0.$$

where $G_r = G_1 \cup G_2$

Rule: When $\beta = -\frac{1}{4}$, $D_{rk}^2 = \mathcal{L}^2$



iii) Average linkage: $D_{pq} = \frac{1}{n_p n_q} \sum_{i \in G_p} \sum_{j \in G_q} d_{ij}$.

Proof: If $G_r = G_p \cup G_q$, $n_r = n_p + n_q$. Then we have:

$$D_{rk}^2 = \frac{n_p}{n_r} D_{pk}^2 + \frac{n_q}{n_r} D_{qk}^2, \quad k \neq p, q.$$

iv) Centroid Method: $D_{rk} = (\bar{X}_k - \bar{X}_r)^T (\bar{X}_k - \bar{X}_r)$

Proof: If $G_r = G_p \cup G_q$, $n_r = n_p + n_q$. Then $\bar{X}_r =$

$$\frac{1}{n_r} (n_p \bar{X}_p + n_q \bar{X}_q). \text{ It's easy to check:}$$

$$D_{rk}^2 = \frac{n_p}{n_r} D_{pk}^2 + \frac{n_q}{n_r} D_{qk}^2 - \frac{n_p n_q}{n_r^2} D_{pq}^2, \quad k \neq p, q.$$

v) Ward Method: It based on ANOVA. $D_{pk}^2 = W_r - W_p - W_q$.

where $G_r = G_p \cup G_q$, $W_p = \sum_i^{n_p} (X_{(i)}^{(p)} - \bar{X}^{(p)})^T (X_{(i)}^{(p)} - \bar{X}^{(p)})$

Proof: It merges two group whose distance is min. i.e.

\hat{Var} of $G_r = G_p \cup G_q$ is min.

Note $\bar{X}_r = \frac{n_p}{n_r} \bar{X}_p + \frac{n_q}{n_r} \bar{X}_q \Rightarrow D_{rk}^2 = \frac{n_p n_q}{n_r} (\bar{X}_p - \bar{X}_q)^T (\bar{X}_p - \bar{X}_q)$

It differs from iv) a constant.

If $G_r = G_p \cup G_q$. Then $D_{rk}^2 = \frac{n_k + n_p}{n_r + n_k} D_{pk}^2 + \frac{n_k + n_q}{n_r + n_k} D_{qk}^2 - \frac{n_k}{n_r + n_k} D_{pq}^2$

⊙ Properties:

i) Monotone:

Def: If in Hierarchical cluster process, D_i is min distance in i^{th} stage, st. $D_1 \leq D_2 \leq \dots \leq D_k$. Then we say the distance we used is monotone.

prop. The distance ① i). iii). v) are monotone but not ① ii) (when $\beta = \frac{1}{4}$). iv)

pf: 1') It's easy to check.

2') $d_{AB}^2 = d_{AC}^2 = 1.1$. $d_{BC}^2 = 1$. for points A, B, C.

is counterexample for the latter.

ii) At each k^{th} stage. $D_{ij}^L \geq D_{ij}^S$. $\forall i, j$.

③ Classes:

i) Def (ρ_{∞}): $T > 0$ is fixed. If G satisfies $d_{ij} \leq T$ for $\forall i, j \in G$. Then G is class w.r.t T .

Rmk: Or we can define: $\frac{1}{|G|-1} \sum_{i \neq j \in G} d_{ij} \leq T$. $\forall i \in G$.

Or $\frac{1}{n(n-1)} \sum_{i \neq j \in G} d_{ij} \leq T$. $d_{ij} \leq M$ $M > T$. $n = |G|$. $\forall i, j \in G$.

Or any partition of $G = G_1 + G_2$. $D(G_1, G_2) \leq T$.

ii) To determine the number of classes:

• Set a threshold value. e.g. $T=2$. in the example

Then $k=3$.

• Observe the data plot

• By statistics: e.g. If $G = \sum_{i=1}^k G_i$. $T_k = \sum_{i=1}^k \frac{1}{n_i} \sum_{j \in G_i} \|X_{ij} - \bar{X}\|^2$

$$= \sum_{i=1}^k W_i + \sum_{i=1}^k n_i (\bar{X}_i - \bar{X})^T (\bar{X}_i - \bar{X}) = P_k + B_k$$

Define: $R_k^z = B_k / T_k$. When $R_k^z \uparrow$, effort \uparrow .

Rule: Refine: $\bar{R}_k^z = R_{k+1}^z - R_k^z$. $\bar{R}_k^z \uparrow$ means

set $G = \sum_{i=1}^{k+1} G_i$ has more effort.

(3) Progressive Cluster:

Idea: Employ a coarse classification

\Rightarrow make adjustment on it basing on some rule

① Initial Classification:

i) 选取凝聚点:

- Base on information or experience. • Randomly.
- 人为分为 k 类. 计算每一类的重心.
- 确定某 r 为半径. 以每个样品为中心做一个球 $r = r$. 确定 D . 取第一个凝聚点为 X_{i_1} . $B(X_{i_1}, r)$ 内的样品取最多. 再依次选点.
- 它与其它凝聚点的距离 $> D$ 且球内样品数最多.

ii) 依据凝聚点初始分类:

- 将每个样品归进距其最近的凝聚点的类中.
- 对数据做标准化处理. $X_{ij} = \frac{X_{ij} - \min_j X_{ij}}{\max_j X_{ij} - \min_j X_{ij}}$. $l = \max_{ij} X_{ij} - \min_{ij} X_{ij}$.

若分为 k 类. 对样品 $X_i = (X_{i1}, \dots, X_{in})$ 计算:

$$\frac{(k-1)(X_{ij} - \min_i X_{ij})}{k} + 1$$
 取 $l \in \mathbb{N}^+$ 距其最近. 归进 l 类.

② Stepwise Adjustment:

1) 初始: 选定凝聚点 \Rightarrow 按距离最近原则分类.

重新计算每一类的重心, 以之作为新的凝聚点.
重复这个过程. 若某一次选出的凝聚点与上一次的完全一致, 则停止.

2) 多类函数: 若有样本 $\{X_i\}_1^n$ 分为 k 类 $\{G_i\}_1^k$.

标记重心为 $\{\bar{X}_k\}$. $|G_i| = n_i$. t_i 为 X_i 所属类

Def: $d(X_i, G_j) = (X_i - \bar{X}_j)^T (X_i - \bar{X}_j)$

$$L(G_1, \dots, G_k) = \sum_{i=1}^n d(X_i, G_{t_i}) \quad (SSD)$$

Prmk: The process above is to make $L(G_1, \dots, G_k)$ to minimal.

Ex: $\{X_i\}_1^5 = \{1, 4, 5, 7, 11\}$. $d_{ij} = |X_i - X_j|$

取 $k=2$. $D = 4k = 4$.

	X_1	X_2	X_3	X_4	X_5
距离	0	1	2	1	0

1) 第一取 X_3 为凝聚点. 然后再取 X_1, X_5 .

2) 计算重心: $G_1^0 = \{X_2, X_3, X_4\}$, $G_2^0 = \{X_1\}$.

$$G_3^0 = \{X_5\}, \quad \bar{X}_1^0 = \frac{14}{3}, \quad \bar{X}_2^0 = 1, \quad \bar{X}_3^0 = 11$$

3) 重新按凝聚点 $\{\frac{14}{3}, 1, 11\}$ 就近分类. 新的

凝聚点仍为 $\frac{14}{3}, 1, 11$. 过程终止.

③ Pointwise Adjustment: (k-method)

一个自然的想法是不用一次性修改一批凝聚点
而是对每个样品归类同时改变凝聚点。

- 1) 确定初始凝聚点数 k ，取前 k 个样品为 k 个凝聚点。
- 2) 确定距离函数，计算类间最小距离 c 与类内最大距离 r 。
- 3) 计算凝聚点两两距离，若 $< c$ ，则合并取重心为新的凝聚点，直到所有凝聚点间距 $> c$ ，此时数量 $k_0 \leq k$ 。
- 4) 对剩下 $n - k_0$ 个样品逐一归类：若该样品与凝聚点最小距 $> r$ ，则单独作为凝聚点，否则 $\leq r$ ，则与最近的凝聚点合并，重新计算重心，作为新凝聚点， \Rightarrow 重复 3)
- 5) 再将样品 $(n - k_0)$ 依次归类按 4) 分类，若某样品二次归类结果一致，则不必计算重心。

例 7. $\{X_i\}_1^5 = \{1, 4, 5, 7, 11\}$, $k=3$, $c=2$, $r=3$, 欧式距离。

1) $d(x_1, x_2) = 3$, $d(x_1, x_3) = 4$, $d(x_2, x_3) = 1 < c = 2$,

合并 x_2, x_3 ，凝聚点为 1, $\frac{4+5}{2}$ ，且已有

两类 $G_1 = \{x_1\}$, $G_2 = \{x_2, x_3\}$, $d(x_1, G_2) = 3.5 > 2$ 。

2) 对剩下的 x_4, x_5 分类, $d(G_1, x_4) = 6$, $d(G_2, x_4)$

$= 2.5 < 3$, x_4 分入 G_2 ，重心为 $\frac{16}{3}$ ；而 $d(G_1, x_5)$

$=10$. $\mu(G_1, X_5) = \frac{17}{3} > 3$. 因此 X_5 单独作为
凝聚点成为一类.

3) 重新分类, 结果仍为 $G_1 = \{X_1\}$, $G_2 = \{X_2, X_3, X_4\}$, $G_3 = \{X_5\}$.

Remark: 最终分类与样品考虑顺序有关, 若

按 $X_5 \rightarrow X_1$, k, c, k 不变, 则结果为:

$G_1 = \{X_5\}$, $G_2 = \{X_4\}$, $G_3 = \{X_1, X_2\}$, $G_4 = \{X_3\}$.

因此前 k 个凝聚点最好取有代表性点.

(4) Fisher's Algorithm:

在某些样品分类中, 我们要求对 X_1, \dots, X_n
分类后的次序不变, i.e. 每一类为 $\{X_{i_1}, \dots, X_{i_{t_k}}\}$
的形式, 此外若分为 k 类, 则共 $\binom{n}{k-1}$ 种.

Process:

i) 类的直径:

设 $G = \{X_t\}_{t=i}^j$, 则记 $G = \{i, i+1, \dots, j\}$, $\bar{X}_G = \frac{1}{j-i+1} \sum_{t=i}^j X_t$.

$D(i, j) = \sum_{t=i}^j (X_t - \bar{X}_G)^T (X_t - \bar{X}_G)$, 为 G 的直径.

ii) Loss Function:

记分法 $b(n, k)$ 为: G_1, G_2, \dots, G_k , $G_i = \{j_i, \dots, j_{i+1}-1\}$.

loss function: $L(ben, k) = \sum_{i=1}^k D(i, i_{i+1}-1)$, i.e.

当 n, k 固定时, 我们要使 $L(ben, k)$ 达到极小, 此时, 各表内的离差最小, 记 $P(n, k)$ 为 n 表最优分表法。

iii) 递推公式:

$$L(P(n, 2)) = \min_{2 \leq j \leq n} \{ D(1, j-1) + D(j, n) \}.$$

$$L(P(n, k)) = \min_{k \leq j \leq n} \{ L(P(i, k-1)) + D(j, n) \}.$$

iv) 求法:

从 $k=2$ 开始, 计算 $L(P(j, 2))$, $j \in [2, n]$

若目前已知 $L(P(i, k-1))$, $i \in [k-1, n]$.

则可寻找 i_k , 使 $L(P(j, k)) = L(P(i_k, k-1))$

+ $D(i_k, j)$, 得到 $\{L(P(i, j)), 1 \leq i \leq n, i \leq j \leq n\}$.

例 1) 计算 $\{D(i, j)\} \Rightarrow L(P(j, 2))$.

2) 递推出 $\{L(P(i, j))\}_{i, j}$, 作图:

$k \backslash i$	2	3	...
3			
4			
5			
...			

$L(P(i, k))$

v) 表数确定:

作图, 取 k 为拐点:

(取 $k=3$ 或 4).

