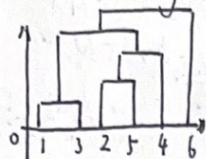# 清华大学

## L12 聚类分析

### 1. Hierachical

Agglomerative 自下而上
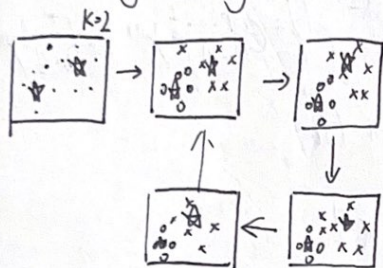
距离
- 欧式
- single (MIN): 对形状特别敏感
- complete (MAX): 不很爱异常值影响
- Group Average
- Ward's linkage (SSE)
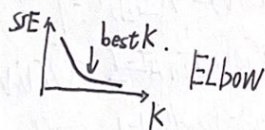


X 方向上坐标无需有序
y 方向上高低一定要格外同.

保证: 什么度越高, 不能出错.

### 2. Partitioning Clustering (K-Means)



Picking K:

① $SSE = \sum_{i=1}^{k} \sum_{x \in C_i} dist^2 (m_i, x)$


best k. ELbow

② Empirical $K \approx \sqrt{\frac{n}{2}}$

③ Cross validation.

## L11 LDA 判别分析

Fisher Approach: Assuption $\Sigma_1 = \Sigma_2$

$\max\limits_{\vec{a}} \dfrac{(\vec{a}^T (\bar{x}_1 - \bar{x}_2))^2}{\vec{a}^T S_{pooled} \vec{a}}$

$\hat{a} = S_{pooled}^{-1} \vec{d}$  $\vec{d} = \bar{x}_1 - \bar{x}_2$

$S_{pooled}^2 = \dfrac{(n_1 - 1) S_1 + (n_2 - 1) S_2}{n_1 + n_2 - 2}$  $y = \hat{a}^T \vec{x}$

maximum $D^2 = (\bar{x}_1 - \bar{x}_2)^T S_{pooled}^{-1} (\bar{x}_1 - \bar{x}_2)$

---

ECM方法: $p_1, p_2$ 为 $\pi_1, \pi_2$ 先验概率 需知 $f_1, f_2$, 先验概率, 以及代价.

$ECM \triangleq C(2|1) P(x \in R_2, x \in \pi_1) + C(1|2) P(x \in R_1, x \in \pi_2)$

$= C(2|1) p_1 + \int_{R_1} [C(1|2) p_2 f_2(x) - C(2|1) p_1 f_1(x)] dx$ 求 min.



$\Rightarrow R_1 = \left\{ x : \dfrac{C(1|2) p_2}{C(2|1) p_1} < \dfrac{f_1(x)}{f_2(x)} \right\}$  $R_2 = R_1^C$

正态分布下的结论

$\Sigma_1 = \Sigma_2$  $f_i(x) = \dfrac{1}{(2\pi)^{p/2} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - u_i)^T \Sigma^{-1} (x - u_i) \right\}$

$R_1: (u_1 - u_2)^T \Sigma^{-1} x - \frac{1}{2} (u_1 - u_2)^T \Sigma^{-1} (u_1 + u_2) \geq \ln \left[ \left( \dfrac{C(1|2)}{C(2|1)} \right) \dfrac{p_2}{p_1} \right]$

sample $\Rightarrow (u_1 - u_2)^T S_{pooled}^{-1} x - \frac{1}{2} (\bar{x}_1 - \bar{x}_2) S_{pooled}^{-1} (\bar{x}_1 + \bar{x}_2) \geq \ln \left[ \dfrac{C(1|2)}{C(2|1)} \dfrac{p_2}{p_1} \right]$

if $\dfrac{C(1|2)}{C(2|1)} \cdot \dfrac{p_2}{p_1} = 1$  $\Rightarrow \hat{y} \geq \frac{1}{2} (\bar{y}_1 + \bar{y}_2) \Rightarrow$ LDA

$\Sigma_1 \neq \Sigma_2$  QDA.

$R_1: -\frac{1}{2} x^T (\Sigma_1^{-1} - \Sigma_2^{-1}) x + (u_1^T \Sigma_1^{-1} - u_2^T \Sigma_2^{-1}) x - k \geq \ln \left[ \dfrac{b(1|2)}{C(2|1)} \dfrac{p_2}{p_1} \right]$

APER: Confusion matrix

|  | | predict membership | | |
|---|---|---|---|---|
|  |  | $\pi_1$ | $\pi_2$ | |
| Actual | $\pi_1$ | $n_{1C}$ | $n_{1m}$ | $n_1$ |
| membership | $\pi_2$ | $n_{2m}$ | $n_{2C}$ | $n_2$ |

$APER = \dfrac{n_{1m} + n_{2m}}{n_1 + n_2}$

Cross validation. 样本均分为n份, 每份中用其他样本作为训练集, 这份验证.

## L8~9 因子分析

$\psi = \begin{bmatrix} \psi_1 & \cdots & \\ & & \psi_p \end{bmatrix}$

$X - u = \underset{p \times m}{L} \underset{m \times 1}{F} + \underset{p \times 1}{\Sigma}$

$E(F) = 0$  $E(\Sigma) = 0$  $Cov\, F = I_{m \times m}$  $Cov(\Sigma) = \underset{p \times p}{\psi}$

$Cov(\Sigma, F) = \underset{p \times m}{O}$

PCA法

$\underset{p \times m}{S} = \hat{L} \hat{L}^T + \underset{p \times p}{\hat{\psi}}$  $\hat{L} = [\sqrt{\lambda_1} e_1, \cdots, \sqrt{\lambda_m} e_m]$

special variance: $S - \hat{L} \hat{L}^T$

$l_{i1}^2 + \cdots + l_{pi}^2 = \hat{\lambda}_i$

proportion $\dfrac{\hat{\lambda}_1 + \cdots \hat{\lambda}_m}{tr(S)}$

MLE法:

$$L(u,\Sigma)=(2\pi)^{-\frac{(n-1)p}{2}}|\Sigma|^{-\frac{(n-1)}{2}}\exp\left\{-\frac{1}{2}tr\left[\Sigma^{-1}\left(\sum_{j=1}^{n}(x_j-\bar{x})(x_j-\bar{x})^T\right)\right]\right\}$$

$$+tr\overline{\left[-\frac{n}{2}\Sigma^{-1}(\bar{x}-u)\right]}\times(2\pi)^{-\frac{p}{2}}|\Sigma|^{-\frac{1}{2}}\exp\left\{-\frac{n}{2}(\bar{x}-u)^T\Sigma^{-1}(\bar{x}-u)\right\}$$
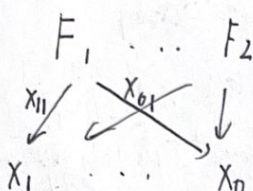
假设: common factors 和 specific factors 都是正态的

$$F\sim N_m(0,I)\quad \varepsilon\sim N_p(0,\psi)\quad F\perp\varepsilon.$$

commonalities:

$$\hat{h}_i^2=\hat{\ell}_{i1}^2+\cdots+\hat{\ell}_{im}^2\quad for\ i=1,2,\cdots p\quad \text{共性方差}$$

$$\psi_i=\sigma_{ii}-\hat{h}_i^2\quad \text{特殊方差}.$$

$$\Sigma=\begin{bmatrix}\sigma_{ii}\\ \square\end{bmatrix}=\begin{bmatrix}\ell_i^T\\ \square\ \ell_i\end{bmatrix}+\begin{bmatrix}\psi_i\\ \square\end{bmatrix}$$

$$F_1\quad\cdots\quad F_2$$



$$X_1\quad\cdots\quad X_0$$

$$\left(\begin{array}{l}\text{Proportion of total sample}\\ \text{variance due to } j\text{-th factor}\end{array}\right)=\frac{\hat{\ell}_{1j}^2+\hat{\ell}_{2j}^2+\cdots+\hat{\ell}_{pj}^2}{S_{11}+\cdots+S_{pp}}\quad \begin{array}{l}\text{第 } j\text{ 向}\\ \text{解释比例}\end{array}$$

$$\hat{h}_i^2=\hat{\ell}_{i1}^2+\hat{\ell}_{i2}^2+\cdots+\hat{\ell}_{im}^2$$

$$\Sigma=E\left((X-u)(X-u)^T\right)=E\left((LF+\varepsilon)(LF+\varepsilon)^T\right)$$

$$=E\left(LFF^TL^T+\varepsilon F^TL^T+LF\varepsilon^T+\varepsilon\varepsilon^T\right)$$

$$=LE(FF^T)L^T+E(\varepsilon\varepsilon^T)=LIL^T+\psi=LL^T+\psi$$

$$\hat{L}^*=V^{\frac{1}{2}}\hat{L}_Z\quad \hat{\psi}^*=V^{\frac{1}{2}}\psi_Z V^{\frac{1}{2}}\quad V=diag(\Sigma)=diag(\sigma_{ii})$$

V 即正规化项.

两种做法优劣: PC: 随 m 增大计算方便

MLE: 相关数与协方差之间系亚可直接转化.

选商更合适的:看可解释性.

因子个数限制 m<p.

Fisher Approach

目标: $\dfrac{\bar{y}_1-\bar{y}_2}{S_y}\quad S_y^2=\dfrac{\sum(y_{1j}-\bar{y}_1)^2+\sum(y_{2j}-\bar{y}_2)^2}{n_1+n_2-2}$

优化: $\max\limits_{a}\dfrac{(\bar{y}_1-\bar{y}_2)^2}{S_y^2}$

$\hat{y}=\hat{a}^T X=(\bar{x}_1-\bar{x}_2)^T S_{pooled}^{-1} X\quad s.t.\ \max\limits_{a}\dfrac{(a^T(\bar{x}_1-\bar{x}_2))^2}{a^T S_{pooled}\, a}$

典型相关分析:

$$X^{(1)}=\begin{bmatrix}X_1^{(1)}\\ \vdots\\ X_p^{(1)}\end{bmatrix}\quad X^{(2)}=\begin{bmatrix}X_1^{(2)}\\ \vdots\\ X_q^{(2)}\end{bmatrix}\quad \Sigma=\begin{bmatrix}\Sigma_{11}&\Sigma_{12}\\ \Sigma_{21}&\Sigma_{22}\end{bmatrix}_{9\times9}$$

目标: 找到最佳 $a,b$ st. $corr(U,V)$ 最大 $U=a^TX^{(1)}$, $V=b^TX^{(2)}$

$$\Rightarrow\max\limits_{a,b\neq0}\frac{a^T\Sigma_{12}b}{\sqrt{a^T\Sigma_{11}a}\sqrt{b^T\Sigma_{22}b}}=\max\frac{\tilde{a}^T\Sigma_{11}^{-\frac{1}{2}}\Sigma_{12}\Sigma_{22}^{-\frac{1}{2}}\tilde{b}}{\sqrt{\tilde{a}^T\tilde{a}}\sqrt{\tilde{b}^T\tilde{b}}}=\max\frac{\tilde{a}^T\Sigma_{11}^{-\frac{1}{2}}\Sigma_{12}\Sigma_{22}^{-\frac{1}{2}}\tilde{b}}{|\tilde{a}|=|\tilde{b}|=1}$$

$\tilde{a}=\Sigma_{11}^{\frac{1}{2}}a\quad \tilde{b}=\Sigma_{22}^{\frac{1}{2}}b$

不要忘记开方.

对 $\Sigma_{11}^{-\frac{1}{2}}\Sigma_{12}\Sigma_{22}^{-\frac{1}{2}}$ 作奇异值分解 $\Rightarrow$ 求 $\Sigma_{11}^{-\frac{1}{2}}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11}^{-\frac{1}{2}}$ 的特征值

假第 $i$th 特征值为 $\rho_i^{*2}$ $\quad U=e_i^T\Sigma_{11}^{-\frac{1}{2}}X^{(1)}\quad V=f_i^T\Sigma_{22}^{-\frac{1}{2}}X^{(2)}$

$$f_i=\frac{1}{\rho_i^*}\Sigma_{22}^{-\frac{1}{2}}\Sigma_{21}\Sigma_{11}^{-\frac{1}{2}}\vec{e}_i$$

$U=AX^{(1)}\quad V=BX^{(2)}$

$$Cov(U,V)=A\Sigma_{12}B^T=\begin{bmatrix}\rho_1^*\\ &\ddots\\ &&\rho_p^*\end{bmatrix}$$

$Cov(U,U)=A\Sigma_{11}A^T=I\quad Cov(V,V)=B\Sigma_{12}B^T=I$

$Var(U_k)=Var(V_k)=1\quad Cov(U_k,V_k)=corr(U_k,V_k)=0\ k\neq\ell$

$Cov(V_k,V_\ell)=corr(V_k,V_\ell)=0\ k\neq\ell$

$Corr(U,X^{(1)})=A\Sigma_{11}V_{11}^{-\frac{1}{2}}$

$Corr(U,X^{(2)})=A\Sigma_{12}V_{22}^{-\frac{1}{2}}$

$Corr(V,X^{(2)})=B\Sigma_{22}V_{22}^{-\frac{1}{2}}\Leftarrow$

$Corr(V,X^{(1)})=B\Sigma_{21}V_{11}^{-\frac{1}{2}}$

$Cov(U,X^{(1)})=Cov(AX^{(1)},X^{(1)})=A\Sigma_{11}$

$Corr(U,X^{(1)})=Cov(U,V_{11}^{-\frac{1}{2}}X^{(1)})$
$=A\Sigma_{11}V_{11}^{-\frac{1}{2}}$

$V_{11}=diag(\Sigma_{11})$

标准化不会影响结果

从几何上来看, 标准化后 $X^{(1)}$ $X^{(2)}$ 张成空间不变, $\therefore U, V$ 也不变.