



聚类分析

多元量个体间距离

两个

例子：计算向量(0,0)、(1,0)、(0,2)两两间的标准化欧氏距离
(假设两个分量的标准差分别为0.5和1)。

(1) 欧氏距离

$$d = \sqrt{\sum (x_i - y_i)^2}$$

$$\sqrt{(1-0)^2 + (0-0)^2} = 1, \quad \sqrt{(0-0)^2 + (2-0)^2} = 2$$

$$\sqrt{(1-0)^2 + (0-2)^2} = \sqrt{5}$$

(2) 曼哈顿距离

$$d = \sum |x_i - y_i|$$

$$|0-1| + |0-0| = 1; \quad |0-0| + |0-2| = 2; \quad |1-0| + |0-2| = 3$$

(3) 切比雪夫距离

$$d = \max |x_i - y_i|$$

$$1, 2, 2$$

(4) 明氏距离

$$d = \sqrt[p]{\sum |x_i - y_i|^p}$$

$p=1 \rightarrow$ 曼哈顿
 $p=2 \rightarrow$ 欧氏
 $p \rightarrow +\infty \rightarrow$ 切比雪夫

(5) 标准化欧氏距离

$$d = \sqrt{\sum \left(\frac{x_{ik} - x_{jk}}{s_k} \right)^2}$$

$$\sqrt{\left(\frac{0-1}{0.5} \right)^2 + \left(\frac{0-0}{1} \right)^2} = 2$$

$$\sqrt{\left(\frac{0-0}{0.5} \right)^2 + \left(\frac{0-2}{1} \right)^2} = 2$$

$$\sqrt{\left(\frac{1-0}{0.5} \right)^2 + \left(\frac{0-2}{1} \right)^2} = 2\sqrt{2}$$

(6) 马氏距离

$$d(x_i, x_j) = \sqrt{(x_i - x_j)^T S^{-1} (x_i - x_j)}$$

已知二维正态总体G的分布为： $G \sim N(\mu, \Sigma)$ ，其中

$$\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} 1 & -0.9 \\ -0.9 & 1 \end{pmatrix}$$

分别求点A=(1, 1)^T和点B=(1, -1)^T到均值μ的欧氏距离和马氏距离。

解：点A到μ的欧氏距离= $\sqrt{1^2 + 1^2} = \sqrt{2}$ ，点B到μ的欧氏距离= $\sqrt{1^2 + 1^2} = \sqrt{2}$ 。

$$\Sigma^{-1} = \frac{1}{0.19} \begin{pmatrix} 1 & 0.9 \\ 0.9 & 1 \end{pmatrix}$$

$$\text{点A到}\mu\text{的马氏距离} = \sqrt{\frac{1}{0.19} \begin{pmatrix} 1 & 0.9 \\ 0.9 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix}} = \sqrt{20}$$

$$\text{点B到}\mu\text{的马氏距离} = \sqrt{\frac{1}{0.19} \begin{pmatrix} 1 & 0.9 \\ 0.9 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix}} = \sqrt{1.05}$$

二维及多个体间距离

(1) 简单匹配函数

	个体 Y_0	个体 Y_1
个体 X_0	a个	b个
个体 X_1	c个	d个

$$S(x, y) = \frac{b + c}{a + b + c + d}$$

(2) Jaccard 函数

$$S(x, y) = \frac{b + c}{b + c + d}$$

字符串相似度的相似性度量

(1) 余弦相似度 $\cos\theta = \frac{x^T y}{\|x\| \|y\|}$

如 $X=(3,2,0,5,0,0,2,0,0)^T$ 和 $Y=(1,0,0,0,0,0,1,0,2)^T$, 它们的余

弦相似度为 $\cos\theta = \frac{3+2}{\sqrt{3^2+2^2+5^2+2^2} \cdot \sqrt{1^2+1^2+2^2}} = 0.31$.

(2) 汉明距离: 两个字符串对应位置不同字符的个数

(3) Jaccard 相似性函数

Jaccard 距离

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$$J_S(A, B) = 1 - J(A, B) = \frac{|A \cup B| - |A \cap B|}{|A \cup B|}$$

如集合 $X=\{1,2,3,4\}$; $Y=\{3,4,5,6\}$;

那么 $J(X, Y) = |\{3,4\}| / |\{1,2,3,4,5,6\}| = 1/3$;

例子: 有两个物品 A, B, 调查 7 位用户是否购买了这两样物品,

得以下向量: $A=(0,0,1,1,1,0,1)$, $B=(1,0,1,0,1,0,0)$

$$|A \cap B| = 2$$

$$|A \cup B| = 5$$

向量 A: (0 0 1 1 1 0 1)

向量 A: (0 0 1 1 1 0 1)

向量 B: (1 0 1 0 1 0 0)

向量 B: (1 0 1 0 1 0 0)

注意, 因为忽略 0-0 匹配。所以 $|A \cup B| \neq 7$ 。

因此, AB 的杰卡德距离为 $1 - \frac{2}{5} = 0.6$ 。

(4) Pearson 相关系数

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y} = \frac{E(XY) - E(X)E(Y)}{\sigma_X \sigma_Y}$$

$$r = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2} \sqrt{\sum (Y_i - \bar{Y})^2}}$$

例1: 计算压力 x 和压缩量 y 之间的相关系数 r 。

$$\bar{X} = \frac{1+2+3+4+5}{5} = 3, \quad \bar{Y} = \frac{1+1+2+2+4}{5} = 2$$

表 2-3 绝缘材料的压缩量和压力表

压力 x (10 lb/in ²)	压缩量 y (0.1 in)
1	1
2	1
3	2
4	2
5	4

$$\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = \sum_{i=1}^5 (X_i - 3)(Y_i - 2) = 7;$$

$$\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} = \sqrt{10}, \quad \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2} = \sqrt{6};$$

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}} = \frac{7}{\sqrt{10}\sqrt{6}} = 0.904$$

类间距离聚类

例子: 二维空间中有6个点, 分别是 $x_1, x_2, x_3, x_4, x_5, x_6$,

数据如表所示。用最近(短)距离法对这6个点进行层次聚类。

最近/最远/中间

数据集		
	x	y
x_1	1	1
x_2	2	1
x_3	1	3
x_4	4	1
x_5	4	4
x_6	5	4

Euclidean
Distance

(1) 最短距离

① 计算初始距离

数据集		
	x	y
x_1	1	1
x_2	2	1
x_3	1	3
x_4	4	1
x_5	4	4
x_6	5	4

D_{ij}	C_1	C_2	C_3	C_4	C_5	C_6
C_1	0					
C_2	1	0				
C_3	4	5	0			
C_4	9	4	13	0		
C_5	18	13	10	9	0	
C_6	25	18	17	10	1	0

② 合并及析

$D^1(1)$ $C_7 = \{x_1, x_2\}$ C_3 C_4 $C_8 = \{x_5, x_6\}$

$C_7 = \{x_1, x_2\}$ 0 C_7 与 C_3, C_4 合并成 C_9

C_3 4 0

C_4 4 13 0

$C_8 = \{x_5, x_6\}$ 13 10 9 0

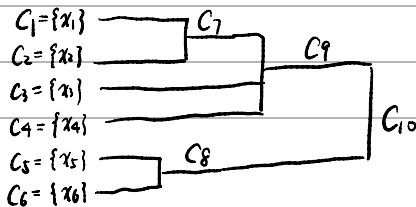
$D^2(2)$ $C_9 = \{x_1, x_2, x_3, x_4\}$ $C_8 = \{x_5, x_6\}$

$C_9 = \{x_1, x_2, x_3, x_4\}$ 0

$C_8 = \{x_5, x_6\}$ 9 0

③ 最后 = 合并

把 C_8, C_9 合成一类 C_{10}



(2) 最大距离

① 十年初始距离

$D^{(0)}$	C_1	C_2	C_3	C_4	C_5	C_6
C_1	0					
C_2	1	0				
C_3	4	5	0			
C_4	9	4	13	0		
C_5	18	13	10	9	0	
C_6	25	18	17	10	1	0

C_1, C_2 合并成 C_7
 C_5, C_6 合并成 C_8

② 合并、更新

$$D^{(1)} \quad C_7 = \{x_1, x_2\} \quad C_3 \quad C_4 \quad C_8 = \{x_5, x_6\}$$

$C_7 = \{C_1, C_2\}$	0			
C_3	5	0		
C_4	9	13	0	
$C_8 = \{x_5, x_6\}$	25	17	10	0

C_1, C_2 合并成 C_7

$$D^{(2)} \quad C_9 = \{x_1, x_2, x_3\} \quad C_4 \quad C_8$$

$C_9 = \{x_1, x_2, x_3\}$	0		
C_4	13	0	
$C_8 = \{x_5, x_6\}$	25	10	0

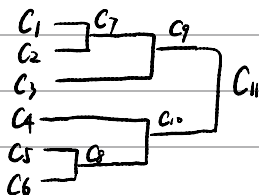
C_4, C_8 合并成 C_{10}

$$D^{(3)} \quad C_9 = \{x_1, x_2, x_3\} \quad C_{10} = \{x_4, x_5, x_6\}$$

$C_9 = \{x_1, x_2, x_3\}$	0	
$C_{10} = \{x_4, x_5, x_6\}$	25	0

③ 排序、= 3 -

C_9, C_{10} 合并成 C_{11}



(3) 中间距离

① 计算初始化距离

D^0	C_1	C_2	C_3	C_4	C_5	C_6
C_1	0					
C_2	<u>1</u>	0				
C_3	4	5	0			
C_4	9	4	13	0		
C_5	18	13	10	9	0	
C_6	25	18	17	10	<u>1</u>	0

② 合并, 更新

$$D^1: C_7 = \{x_1, x_2\} \quad C_3 \quad C_4 \quad C_8 = \{x_5, x_6\}$$

$C_7 = \{x_1, x_2\}$	0			
C_3	<u>4.25</u>	0		
C_4	6.25	13	0	
$C_8 = \{x_5, x_6\}$	18	13.25	9.25	0

$$D_{37}^2 = \frac{1}{2} (D_{13}^2 + D_{23}^2) - \frac{1}{4} D_{12}^2 = \frac{1}{2} (4+5) - \frac{1}{4} \times 1 = 4.25$$

$$D_{47}^2 = \frac{1}{2} (D_{14}^2 + D_{24}^2) - \frac{1}{4} D_{12}^2 = \frac{1}{2} (9+4) - \frac{1}{4} \times 1 = 6.25$$

$$D_{38}^2 = \frac{1}{2} (D_{35}^2 + D_{36}^2) - \frac{1}{4} D_{56}^2 = \frac{1}{2} (10+17) - \frac{1}{4} \times 1 = 13.25$$

$$D_{48}^2 = \frac{1}{2} (D_{45}^2 + D_{46}^2) - \frac{1}{4} D_{56}^2 = \frac{1}{2} (9+10) - \frac{1}{4} \times 1 = 9.25$$

$$D_{78}^2 = \frac{1}{4} (D_{15}^2 + D_{25}^2 + D_{16}^2 + D_{26}^2 - D_{12}^2 - D_{56}^2) = \frac{1}{4} (18+25+13+18-2) = 18$$

$$D^2: C_9 = \{x_1, x_2, x_3\} \quad C_4 = \{x_4\} \quad C_8 = \{x_5, x_6\}$$

$C_9 = \{x_1, x_2, x_3\}$	0		
$C_4 = \{x_4\}$	<u>8.5625</u>	0	
$C_8 = \{x_5, x_6\}$	14.5625	9.25	0

$$D_{49}^2 = \frac{1}{2} (D_{47}^2 + D_{34}^2) - \frac{1}{4} D_{37}^2 = \frac{1}{2} \times (6.25 + 13) - \frac{1}{4} \times 4.25 = 8.5625$$

$$D_{89}^2 = \frac{1}{2} (D_{78}^2 + D_{38}^2) - \frac{1}{4} D_{37}^2 = \frac{1}{2} \times (18 + 13.25) - \frac{1}{4} \times 4.25 = 14.5625$$

$D^2(3)$

$C_{10} = \{x_1, x_2, x_3, x_4\} \quad C_8 = \{x_5, x_6\}$

$C_{10} = \{ \quad \} \quad 0$

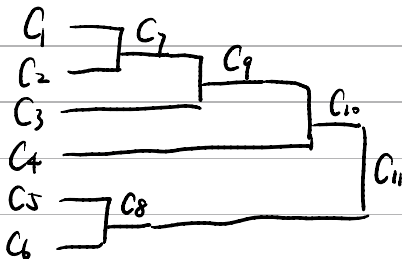
$C_8 = \{ \quad \} \quad 9.765625 \quad 0$

$$D_{810}^2 = \frac{1}{2}(D_{48}^2 + D_{89}^2) - \frac{1}{4}D_{29}^2 = \frac{1}{2} \times (9.25 + 14.5625) - \frac{1}{4} \times 8.5625$$

$$= 9.765625$$

③ 取 $\beta = \frac{1}{2}$

把 C_8, C_{10} 合成一个类 C_{11}



K-Means

假设给定如下要进行聚类的元组: {2,4,10,12,3,20,30,11,25},并

假设 $k=2$ 。

► 初始时用前两个数值作为类的均值; $m_1=2$ 和 $m_2=4$ 。

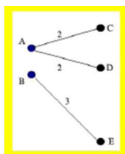
簇1均值	簇2均值	析簇	析簇1均值	析簇2	析簇2均值
2	4	{2,3}	2.5	{4,10,11,12,20,25,30}	16
2.5	16	{2,3,4}	3	{10,11,12,20,25,30}	18
3	18	{2,3,4,10}	4.75	{11,12,20,25,30}	19.6
4.75	19.6	{2,3,4,10,11,12}	7	{20,25,30}	25
4.75	19.6	{2,3,4,10,11,12}	7	{20,25,30}	25

算法已收敛

K-Medoids 算法

假设空间中的五个点{A、B、C、D、E}, 如下图所示, 各点之间的距离关系如下表所示。根据所给的数据对其运行k-medoids算法实现划分聚类(设 $k=2$)。

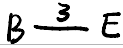
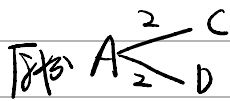
样本点	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0



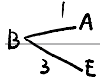
第一步 建立阶段: 假如从5个对象中随机抽取的2个中心点为{A, B}, 则样本被划分为{A, C, D}和{B, E}, 如图所示。

第二步 交换阶段: 假定中心点A、B分别被非中心点{C、D、E}替换, 根据PAM算法需要计算下列代价

TC_{AC} 、 TC_{AD} 、 TC_{AE} 、 TC_{BC} 、 TC_{BD} 、 TC_{BE} 。



$$TC_{AC} = CA_{AC} + CB_{AC} + CC_{AC} + CD_{AC} + CE_{AC}$$

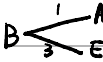


$$= (1-0) + 0 + (0-2) + (1-2) + (3-3)$$

$$= -2$$

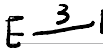


$$TC_{AD} = CA_{AD} + CB_{AD} + CC_{AD} + CD_{AD} + CE_{AD}$$

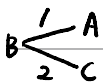


$$= (1-0) + 0 + 0 + (0-2) + (2-3)$$

$$= -2$$

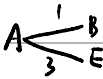


$$TC_{AE} = CA_{AE} + CB_{AE} + CC_{AE} + CD_{AE} + CE_{AE}$$



$$= (1-0) + 0 + 0 + (3-2) + (0-3)$$

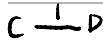
$$= -1$$



$$TC_{BC} = CA_{BC} + CB_{BC} + CC_{BC} + CD_{BC} + CE_{BC}$$

$$= 0 + (1-0) + (0-2) + (1-2) + (3-3)$$

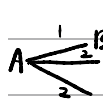
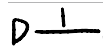
$$= -2$$



$$TC_{BD} = CA_{BD} + CB_{BD} + CC_{BD} + CD_{BD} + CE_{BD}$$

$$= 0 + (1-0) + (1-2) + (0-2) + (3-0)$$

$$= -2$$



$$TC_{BE} = CA_{BE} + CB_{BE} + CC_{BE} + CD_{BE} + CE_{BE}$$

$$= 0 + (1-0) + (2-2) + (2-2) + (0-3)$$

$$= -2$$

E

DIANA算法

例子：有如下表所示的数据集，使用DIANA算法对该数据集进行分裂层次聚类。

序号	属性 1	属性 2
1	1	1
2	1	2
3	2	1
4	2	2
5	3	4
6	3	5
7	4	4
8	4	5

1	0							
2	1	0						
3	1	1.4	0					
4	1.4	1	1	0				
5	3.6	2.8	3.2	2.2	0			
6	4.5	3.6	4.1	3.2	1	0		
7	4.2	3.6	3.6	2.8	1	1.4	0	
8	5	4.2	4.5	3.6	1.4	1	1	0
1	2	3	4	5	6	7	8	

对于所给的数据进行DIANA算法，(设 $n=8$,用户输入的终止条件为2个类)，初始类{1,2,3,4,5,6,7,8}。

① 计算平均距离 (① 找直径最大的类)

序号1的平均距离 (就是1距离其它各个点的距离长度之和除以7)

$$s_1 = (1+1+1.1414+3.6+4.47+4.24+5)/7 = 2.96;$$

$$\text{序列2的平均距离 } s_2 = (1+1.414+1+2.828+3.6+3.6+4.24)/7 = 2.526;$$

$$\text{序列3的平均距离 } s_3 = (1+1.414+1+3.16+4.12+3.6+4.27)/7 = 2.68;$$

$$\text{序列4的平均距离 } s_4 = (1.414+1+1+2.24+3.16+2.828+3.6)/7 = 2.18$$

$$\text{序列5的平均距离 } s_5 = 2.18;$$

$$\text{序列6的平均距离 } s_6 = 2.68;$$

$$\text{序列7的平均距离 } s_7 = 2.526;$$

$$\text{序列8的平均距离 } s_8 = 2.96;$$

② 找平均距离最大的放到 splinter group, 剩下的放 old party

$$\text{splinter group} = \{1\} \quad \text{old party} = \{2, 3, 4, 5, 6, 7, 8\}$$

③ 找离 splinter group 最近的, 若该距离 \leq 该类 old party 中其他点的最小距离, 则舍弃

一直找直到找不到

↓ 1	0									
2	0									
1	1.4	0								
1.4	1	0								
3.6	2.8	3.2	2.2	0						
4.5	3.6	4.1	3.2	1	0					
4.2	3.6	3.6	2.8	1	1.4	0				
5	4.2	4.5	3.6	1.4	1	1	0			

splinter group = {1} {1, 2}

↓ 1	0									
2	1	0								
3	1.4	0								
3	1.4	1	0							
3.6	2.8	3.2	2.2	0						
4.5	3.6	4.1	3.2	1	0					
4.2	3.6	3.6	2.8	1	1.4	0				
5	4.2	4.5	3.6	1.4	1	1	0			

{1, 2, 3}

↓ 1	0									
2	1	0								
3	1	1.4	0							
4	1.4	1	0							
4	1.4	1	0							
3.6	2.8	3.2	2.2	0						
4.5	3.6	4.1	3.2	1	0					
4.2	3.6	3.6	2.8	1	1.4	0				
5	4.2	4.5	3.6	1.4	1	1	0			

找到了

④ splinter group 为一类, old party 为一类

$$\{1, 2, 3, 4\}$$

$$\{5, 6, 7, 8\}$$

分类分析

协方差矩阵相等的距离判别法

例1: 记二维正态总体 $N_2(\mu^{(i)}, \Sigma)$ 为 $G_i (i=1, 2)$ (两总体协方差阵相同), 已知来自 $G_i (i=1, 2)$ 的样本数据为

$$X^{(1)} = \begin{pmatrix} 2 & 12 \\ 4 & 10 \\ 3 & 8 \\ 3 & 10 \end{pmatrix}, X^{(2)} = \begin{pmatrix} 5 & 7 \\ 3 & 9 \\ 4 & 5 \end{pmatrix} \quad (k=2, m=2) \\ (n_1=4, n_2=3)$$

(1) 试求两总体的样本离差阵 S_1, S_2 和合并样本协差阵 S .

(2) 今有样本 $x_0 = (2, 8)'$, 试问按马氏距离准则样本 x_0 应判归哪一类。

样本离差阵

$$S_1 = \sum_{j=1}^{n_1} (x_j^{(1)} - \bar{x}^{(1)})' (x_j^{(1)} - \bar{x}^{(1)})$$

$$S_2 = \sum_{j=1}^{n_2} (x_j^{(2)} - \bar{x}^{(2)})' (x_j^{(2)} - \bar{x}^{(2)})$$

样本合并组内离差阵 $S_1 + S_2$

$$\text{合并样本协差阵: } \frac{1}{n_1+n_2-2} (S_1 + S_2)$$

$$\text{判别函数 } W(x) = (x - \bar{x})' S^{-1} (\bar{x}^{(1)} - \bar{x}^{(2)})$$

(1) ① 求每类的样本均值向量

$$\bar{X}^{(1)}: \frac{2+4+3+3}{4} = 3 \quad \bar{X}^{(1)} = \begin{pmatrix} 3 \\ 10 \end{pmatrix}$$

$$\frac{12+10+8+10}{4} = 10$$

$$\bar{X}^{(2)}: \frac{5+3+4}{3} = 4 \quad \bar{X}^{(2)} = \begin{pmatrix} 4 \\ 7 \end{pmatrix}$$

$$\frac{7+9+5}{3} = 7$$

② 样本均值作差

$$\tilde{X}^{(1)} = \begin{pmatrix} 2-3 & 12-10 \\ 4-3 & 10-10 \\ 3-3 & 8-10 \\ 3-3 & 10-10 \end{pmatrix} = \begin{pmatrix} -1 & 2 \\ 1 & 0 \\ 0 & -2 \\ 0 & 0 \end{pmatrix}$$

$$\tilde{X}^{(2)} = \begin{pmatrix} 5-4 & 7-7 \\ 3-4 & 9-7 \\ 4-4 & 5-7 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ -1 & 2 \\ 0 & -2 \end{pmatrix}$$

③ 自己乘积求样本离差阵

$$S_1 = (\tilde{X}^{(1)})' (\tilde{X}^{(1)}) = \begin{pmatrix} -1 & 1 & 0 & 0 \\ 2 & 0 & -2 & 0 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \\ 0 \\ -2 \end{pmatrix} = \begin{pmatrix} 2 & -2 \\ -2 & 8 \end{pmatrix}$$

$$S_2 = (\tilde{X}^{(2)})' (\tilde{X}^{(2)}) = \begin{pmatrix} 1 & -1 & 0 \\ 0 & 2 & -2 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \\ 0 \\ -2 \end{pmatrix} = \begin{pmatrix} 2 & -2 \\ -2 & 8 \end{pmatrix}$$

④ 相加, 乘系数, 求合并样本协差阵

$$S = \frac{1}{n_1+n_2-2} (S_1 + S_2) = \frac{1}{5} \begin{pmatrix} 4 & -4 \\ -4 & 16 \end{pmatrix} = \frac{4}{5} \begin{pmatrix} 1 & -1 \\ -1 & 4 \end{pmatrix}$$

合并样本协差阵 合并组内离差阵

(2) f_1 : 贝叶斯判别距离

$$\text{先求 } S^{-1} \quad |S| = \frac{16}{25} \begin{vmatrix} 1 & -1 \\ -1 & 4 \end{vmatrix} = \frac{16}{25} \times 3 = \frac{48}{25}$$

$$S^* = \frac{4}{5} \begin{pmatrix} 4 & 1 \\ 1 & 1 \end{pmatrix}$$

$$S^{-1} = \frac{1}{|S|} S^* = \frac{25}{48} \times \frac{4}{5} \begin{pmatrix} 4 & 1 \\ 1 & 1 \end{pmatrix} = \frac{5}{12} \begin{pmatrix} 4 & 1 \\ 1 & 1 \end{pmatrix}$$

$$D_1^2(x_0) = (x_0 - \bar{x}^{(1)})' S^{-1} (x_0 - \bar{x}^{(1)})$$

$$= (-1 \ -2) \frac{5}{12} \begin{pmatrix} 4 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} -1 \\ -2 \end{pmatrix} = \frac{5}{12} (-6 \ -3) \begin{pmatrix} -1 \\ -2 \end{pmatrix} = 5$$

$$D_2^2(x_0) = (x_0 - \bar{x}^{(2)})' S^{-1} (x_0 - \bar{x}^{(2)})$$

$$= (-2 \ 1) \frac{5}{12} \begin{pmatrix} 4 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} -2 \\ 1 \end{pmatrix} = \frac{5}{12} (-7 \ -1) \begin{pmatrix} -2 \\ 1 \end{pmatrix} = \frac{5 \times 3}{12} = \frac{65}{12}$$

$$D_1^2(x_0) < D_2^2(x_0) \quad \therefore x_0 \in G_1$$

f_2 用判别函数

$$W(x) = (x - \bar{X})' S^{-1} (\bar{x}^{(1)} - \bar{x}^{(2)})$$

$$\bar{X} = \frac{1}{2} (\bar{x}^{(1)} + \bar{x}^{(2)}) = \begin{pmatrix} 3.5 \\ 8.5 \end{pmatrix}$$

$$W(x_0) = (-1.5 \ -0.5) \frac{5}{12} \begin{pmatrix} 4 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} -1 \\ 3 \end{pmatrix}$$

$$= \frac{5}{12} (-6.5 \ -2) \begin{pmatrix} -1 \\ 3 \end{pmatrix}$$

$$= \frac{5}{12} \times 0.5 = \frac{5}{24} > 0 \quad \therefore x_0 \in G_1$$

协方差矩阵不相等的距离判别法

例子：已知有两个类 G_1 和 G_2 ，分别为设备 A、B 生产的产品。设备 A 生产的产品平均耐磨度 $\mu_1 = 80$ ，精度 $\sigma_1^2 = 0.25$ ；设备 B 的平均耐磨度 $\mu_2 = 75$ ，精度 $\sigma_2^2 = 4$ 。现有一耐磨度为 78 的产品 x ，试判断它为哪一台设备生产的。

f_1 ① 决策 $\mu^* = \frac{\mu_1 \sigma_2 + \mu_2 \sigma_1}{\sigma_1 + \sigma_2}$ $\mu_* = \frac{\mu_1 \sigma_2 - \mu_2 \sigma_1}{\sigma_1 - \sigma_2}$ f_2

$$\mu^* = \frac{80 \times 2 + 75 \times 0.5}{0.5 + 2} = 79$$

$$\mu_* = \frac{80 \times 2 - 75 \times 0.5}{2 - 0.5} = 81.67$$

首先计算马氏距离：

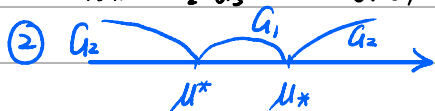
$$D^2(x, G_1) = \frac{(x - \mu_1)^2}{\sigma_1^2} = \frac{(78 - 80)^2}{0.25} = 16 = 4^2$$

$$D^2(x, G_2) = \frac{(x - \mu_2)^2}{\sigma_2^2} = \frac{(78 - 75)^2}{4} = 2.25 = 1.5^2$$

$$D^2(x, G_2) < D^2(x, G_1)$$

$$x = 78 < \mu^*, \quad \therefore x \in G_2$$

$$x \in G_2$$



ID3 算法构造决策树

表 4-1 高尔夫活动决策表

编号	天气	温度	湿度	风速	活动
1.	晴	炎热	高	弱	取消
2.	晴	炎热	高	强	取消
3.	阴	炎热	高	弱	进行
4.	雨	适中	高	弱	进行
5.	雨	寒冷	正常	弱	进行
6.	雨	寒冷	正常	强	取消
7.	阴	寒冷	正常	强	进行
8.	晴	适中	高	弱	取消
9.	晴	寒冷	正常	弱	进行
10.	雨	适中	正常	弱	进行
11.	晴	适中	正常	强	进行
12.	阴	适中	高	强	进行
13.	阴	炎热	正常	弱	进行
14.	雨	适中	高	强	取消

① 信息熵: $Ent(D) = -\sum_{i=1}^m P_i \log_2(P_i)$

② 信息增益: 划分前 - 划分后

$$Gain(D, a) = Ent(D) - \sum_{i=1}^V \frac{|D^i|}{|D|} Ent(D^i)$$

(D^i 表示在特征 a 上取值的所有样本)

③ 增益率: $Gain_ratio(D, a) = \frac{Gain(D, a)}{IV(a)}$

$$IV(a) = -\sum_{i=1}^V \frac{|D^i|}{|D|} \log_2 \frac{|D^i|}{|D|}$$

① 计算样本的熵

$$I(S_1, S_2) = I(9, 5) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.940$$

② 算每个属性的信息增益, 和最大的对树分支, 一直重复直到全叶

天气: 晴天 $I_{晴} = I(S_{11}, S_{21}) = I(2, 3) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$

阴天: $I_{阴} = I(S_{12}, S_{22}) = I(4, 0) = -\frac{4}{4} \log_2 1 = 0$

雨天: $I_{雨} = I(S_{13}, S_{23}) = I(3, 2) = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{2}{7} \log_2 \frac{2}{7} = 0.971$

$$Ent(\text{天气}) = \frac{5}{14} \times 0.971 + \frac{5}{14} \times 0.971 = 0.694$$

$$Gain(\text{天气}) = 0.940 - 0.694 = 0.246$$

温度: 炎热: $I_{炎} = I(S_{11}, S_{21}) = I(2, 2) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$

适中: $I_{中} = I(S_{12}, S_{22}) = I(4, 2) = -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = 0.918$

寒冷: $I_{冷} = I(S_{13}, S_{23}) = I(3, 1) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.811$

$$Ent(\text{温度}) = \frac{4}{14} \times 1 + \frac{4}{14} \times 0.918 + \frac{4}{14} \times 0.811 = 0.911$$

$$Gain(\text{温度}) = 0.940 - 0.911 = 0.029$$

湿度: 高: $I_{高} = I(S_{01}, S_{21}) = I(3, 4) = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.985$

正常: $I_{正} = I(S_{12}, S_{22}) = I(6, 1) = -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} = 0.592$

$$Ent(风速) = \frac{1}{2} \times 0.985 + \frac{1}{2} \times 0.592 = 0.7885$$

$$Gain(风速) = 0.940 - 0.7885 = 0.1515$$

风速: 弱: $I(S_1, S_2) = I(6, 2) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.811$

强: $I(S_1, S_2) = I(3, 3) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$

$$Ent(风速) = \frac{8}{14} \times 0.811 + \frac{6}{14} \times 1 = 0.892$$

$$Gain(风速) = 0.940 - 0.892 = 0.048$$

显然, 天气的信息熵最大, 按天气划分

再对晴中的 $\{1, 2, 8, 9, 11\}$ 划分

$$I(S_1, S_2) = I(2, 3) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

低温: $I(S_{11}, S_{21}) = I(0, 2) = -\frac{2}{2} \log_2 1 = 0$

适中: $I(S_{12}, S_{22}) = I(1, 1) = -\frac{1}{2} \log_2 - \frac{1}{2} \log_2 = 1$

寒冷: $I(S_{13}, S_{23}) = I(1, 0) = \log_2 1 = 0$

$$Gain(低温) = 0.971 - \frac{2}{5} = 0.571$$

高温: $I(S_{14}, S_{24}) = I(0, 3) = -\frac{3}{3} \log_2 1 = 0$

$$I(S_{21}, S_{22}) = I(2, 0) = 0$$

$$Gain(低温) = 0.971 - \frac{2}{5} \text{最大, 不用解风速}$$

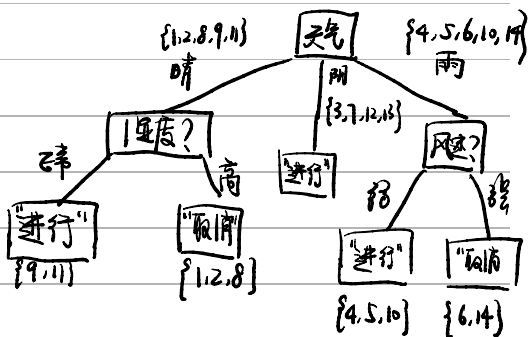
对雨中的 $\{4, 5, 6, 10, 14\}$ 划分

$$I(S_1, S_2) = I(3, 2) = 0.971$$

风速: 弱: $I(S_{14}, S_{21}) = I(3, 0) = 0$

强: $I(S_{12}, S_{22}) = I(0, 2) = 0$

$$Gain(风速) = 0.971 - \frac{2}{5} \text{最大}$$



问题：上面数据集中序号1-12为已知的电影分类，分为喜剧片、动作片、爱情片三个种类，使用的特征值分别为搞笑镜头、打斗镜头、拥抱镜头的数量。那么来了一部新电影《唐人街探案》，它属于上述3个电影分类中的哪个类型？

电影分类数据集（纯属虚构）：

序号	电影名称	搞笑镜头	拥抱镜头	打斗镜头	电影类型
1.	宝贝当家	45	2	9	喜剧片
2.	美人鱼	21	17	5	喜剧片
3.	澳门风云3	54	9	11	喜剧片
4.	功夫熊猫3	39	0	31	喜剧片
5.	谍影重重	5	2	57	动作片
6.	叶问3	3	2	65	动作片
7.	伦敦陷落	2	3	55	动作片
8.	我的特工爷爷	6	4	21	动作片
9.	奔爱	7	46	4	爱情片
10.	夜孔雀	9	39	8	爱情片
11.	代理情人	9	38	2	爱情片
12.	新步步惊心	8	34	17	爱情片
13.	唐人街探案	23	3	17	?

① 算距离(升序排列)

② 前k个里面哪个类就选哪个

1	[['我的特工爷爷', 17.49],	——动作	k=5 认为是喜剧
2	[['美人鱼', 18.55],	——喜剧	
3	[['功夫熊猫3', 21.47],	——喜剧	
4	[['宝贝当家', 23.43],	——喜剧	
5	[['澳门风云3', 32.14],	——喜剧	
6	[['新步步惊心', 34.44],		
7	[['夜孔雀', 39.66],		
8	[['代理情人', 40.57],		
9	[['伦敦陷落', 43.42],		
10	[['谍影重重', 43.87],		
11	[['奔爱', 47.69],		
12	[['叶问3', 52.01]]		

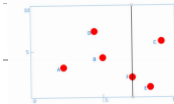
KD树(及其搜索算法)

例子: 给定一个二维空间的数据集: $T = \{(2,3), (5,4), (9,6), (4,7), (8, 1), (7,2)\}$, 构造一个平衡KD树。

根节点对应包含数据集 T 的矩形, 首先选择 $x^{(1)}$ 维, 按照 $x^{(1)}$ 维对数据集进行排序, 得到: $A(2, 3)$ 、 $D(4, 7)$ 、 $B(5, 4)$ 、 $F(7, 2)$ 、 $E(8, 1)$ 、 $C(9, 6)$ 。

对于 $x^{(1)}$ 维, 其中位数为: 7, 选择 $F(7,2)$ 作为根节点。

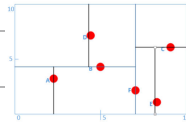
接着按照 $x^{(1)}$ 维进行划分, 将小于7的划分到根节点的左子树中, 大于7的划分到根节点的右子树中(该过程类似搜索二叉树)。对应的树结构如下:



此时对于左右子树, 需要选取相同的维进行划分(此处选取 $x^{(2)}$ 维)。

✓ 左子树: 对于节点 $(2,3), (4,7), (5,4)$, 按照 $x^{(2)}$ 维进行排序 $(2,3), (5,4), (4,7)$, 选择中位数 $(5,4)$ 作为子树根节点; 同理, 将小于4的放到 $B(5,4)$ 节点的左子树中, 大于4的放到 $B(5,4)$ 节点的右子树中。

✓ 对于根节点 $F(7,2)$ 的右子树, 对节点 $(8,1), (9,6)$ 进行排序, 并选择中位数6, 由于 $1 < 6$, 将 $E(8,1)$ 作为 $C(9,6)$ 的左孩子。到此, KD树构造完成。



Naive Bayes

例1: 给出如表所示的训练样本, 目的是判定一个人是否会购买电脑。这个人的属性为 $X = (\text{年龄} < 30, \text{收入} = \text{中等}, \text{学生} = \text{是}, \text{信用} = \text{一般})$ 。

解: 设类别 C_1 : 购买电脑 = “是”, 类别 C_2 : 购买电脑 = “否”, 所以可求:
 $P(C_1) = P(\text{购买电脑} = \text{“是”}) = 9/14 = 0.643$
 $P(C_2) = P(\text{购买电脑} = \text{“否”}) = 5/14 = 0.357$

判定一个人是否会购买电脑的训练样本

编号	年龄	收入	学生	信用等级	类型: 购买电脑
1	<=30	高	否	一般	不会购买
2	<=30	高	否	良好	不会购买
3	31~40	高	否	一般	不会购买
4	>40	中等	否	一般	会购买
5	>40	低	是	一般	会购买
6	>40	低	是	良好	不会购买
7	31~40	低	是	良好	会购买
8	<=30	低	是	一般	不会购买
9	<=30	中	是	一般	不会购买
10	>40	中等	是	一般	会购买
11	<=30	中等	是	良好	会购买
12	31~40	中等	否	良好	会购买
13	31~40	高	是	一般	会购买
14	>40	中等	否	良好	不会购买

计算每个类别的 $P(X|C_i)$:

$P(\text{年龄} = \text{“<30”} | \text{购买电脑} = \text{“是”}) = 2/9 = 0.222$
 $P(\text{年龄} = \text{“<30”} | \text{购买电脑} = \text{“否”}) = 3/5 = 0.6$
 $P(\text{收入} = \text{“中等”} | \text{购买电脑} = \text{“是”}) = 4/9 = 0.444$
 $P(\text{收入} = \text{“中等”} | \text{购买电脑} = \text{“否”}) = 2/5 = 0.4$

$$P(X|C_i) = \prod_{k=1}^n P(X_k | C_i)$$

$P(\text{学生} = \text{“是”} | \text{购买电脑} = \text{“是”}) = 6/9 = 0.667$
 $P(\text{学生} = \text{“是”} | \text{购买电脑} = \text{“否”}) = 1/5 = 0.2$
 $P(\text{信用} = \text{“一般”} | \text{购买电脑} = \text{“是”}) = 6/9 = 0.667$
 $P(\text{信用} = \text{“一般”} | \text{购买电脑} = \text{“否”}) = 2/5 = 0.4$

编号	年龄	收入	学生	信用等级	类型: 购买电脑
1	<=30	高	否	一般	不会购买
2	<=30	高	否	良好	不会购买
3	31~40	高	否	一般	不会购买
4	>40	中等	否	一般	会购买
5	>40	低	是	一般	会购买
6	>40	低	是	良好	不会购买
7	31~40	低	是	良好	会购买
8	<=30	低	是	一般	不会购买
9	<=30	中	是	一般	不会购买
10	>40	中等	是	一般	会购买
11	<=30	中等	是	良好	会购买
12	31~40	中等	否	良好	会购买
13	31~40	高	是	一般	会购买
14	>40	中等	否	良好	不会购买

$$P(X|C_i) = \prod_{k=1}^n P(X_k | C_i)$$

从而得到:
 $P(X|\text{购买电脑} = \text{“是”}) = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$
 $P(X|\text{购买电脑} = \text{“否”}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$

又由于: $P(X|C_i) P(C_i)$

$$P(C_1) = P(\text{购买电脑} = \text{“是”}) = 9/14 = 0.643$$

$$P(C_2) = P(\text{购买电脑} = \text{“否”}) = 5/14 = 0.357$$

所以:

$$P(X|\text{购买电脑} = \text{“是”}) \times P(\text{购买电脑} = \text{“是”}) = 0.028$$

$$P(X|\text{购买电脑} = \text{“否”}) \times P(\text{购买电脑} = \text{“否”}) = 0.007$$

所以判定 X 处于类别 C_1 , 此人会购买电脑。

关联分析

GBDT 树

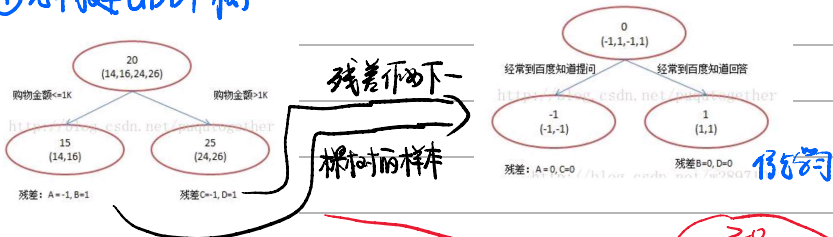
GBDT学习例子：训练集：(A, 14岁), (B, 16岁), (C, 24岁), (D, 26岁)。
训练数据的均值为：20岁；决策树的个数为：2棵。每个样本的特征有两个：购买金额是否小于1K，经常去百度提问还是回答？

测试样本：预测一个购物金额为3K，经常去百度问淘宝相关问题的女生的年龄。

表示如下：

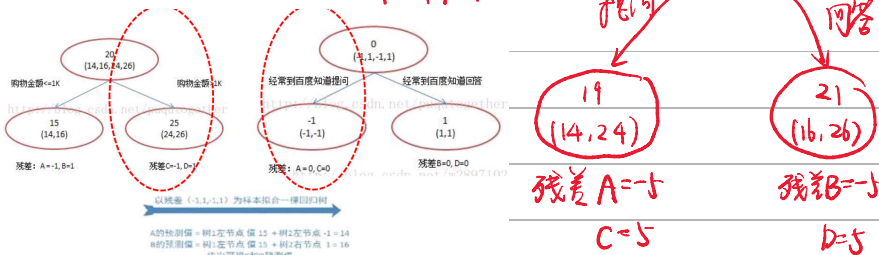
A: 14岁, 购物金额 $\leq 1K$; 经常到百度知道提问
 B: 16岁, $\leq 1K$; 回答
 C: 24岁 $> 1K$; 提问
 D: 26岁 $> 1K$; 回答

①先构建GBDT树



②使用GBDT树进行迭代

为什么不利用第一个模型



- 提取2个特征：购物金额3K，经常去百度上面提问问题。
- 第一棵树→购物金额大于1K→右叶子，初步说明这个女生25岁。
第二棵树→经常去百度提问→左叶子，说明这个女生的残差为-1。
- 叠加前面每棵树得到的结果：25-1=24岁，最终预测结果为24岁。

关联规则分析

① 支持度: 规则的支持度是包含该集的物数与总物数的比值.

S 规则 $X \Rightarrow Y$ 的支持度是 D 中事务同时包含 X, Y 的百分比

② 置信度: 规则 $X \Rightarrow Y$ 的置信度是 D 中已包含 X 的事务中, 包含 Y 的百分比

C

③ 提升度: 置信度 / 含有 Y 的比例

Lift

④ 兴趣因子: $I(X \Rightarrow Y) = \frac{S(A \Rightarrow B)}{S(A) \times S(B)}$

I

= 无交互的提升度 \Leftrightarrow 兴趣因子

⑤ 教育项集: 支持度 > 最低提升度阈值

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example:

{Milk, Diaper} \Rightarrow Beer

$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|T|} = \frac{2}{5} = 0.4$$

Example:

{Milk, Diaper} \Rightarrow Beer

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

关联规则分析

例子: 一个Apriori的具体例子, 该例基于下图的AllElectronics的事务数据库。数据库中有9个事务, 即 $|D|=9$ 。Apriori假定事务中的项按字典次序存放。

TID	商品ID的列表
T100	I1,I2,I5
T200	I2,I4
T300	I2,I3
T400	I1,I2,I4
T500	I1,I3
T600	I2,I3
T700	I1,I3
T800	I1,I2,I3,I5
T900	I1,I2,I3

支持度阈值=2

置信度阈值=70%

求频繁项集

求最高频繁项集产生的关联规则

频繁1项集	候选2项集	2-项集	候选3项集	3-项集
{1} 6	{12} 4	{12} 4	{123} 2	{123} 2
{2} 7	{13} 4	{13} 4	{125} 2	{125} 2
{3} 6	{14} 1	{15} 2	{135} 1	
{4} 2	{15} 2	{23} 4	{234} 0	
{5} 2	{23} 4	{24} 2	{235} 1	
	{24} 2	{25} 2	{245} 0	
	{25} 2			
	{34} 0			
	{35} 1			
	{45} 0			

4-项集 {1235} 不再是频繁项集, 算法结束

求关联规则:

数据子-项集 $\{1, 2, 3\}$ $\{1, 2, 5\}$

先求 $\{1, 2, 3\}$

$\{1, 2, 5\}$

$\{1\}$ $\{2\}$ $\{3\}$ $\{1, 2\}$ $\{1, 3\}$ $\{2, 3\}$

$\{1\}$ $\{2\}$ $\{5\}$ $\{1, 2\}$ $\{1, 5\}$ $\{2, 5\}$

$\{1\} \Rightarrow \{2, 3\}$ $2/6 = 33.3\%$

$\{1\} \Rightarrow \{2, 5\}$ $2/6$

$\{2\} \Rightarrow \{1, 3\}$ $2/7$

$\{2\} \Rightarrow \{1, 5\}$ $2/7$

$\{3\} \Rightarrow \{1, 2\}$ $2/6$

$\{5\} \Rightarrow \{1, 2\}$ $2/2 = 100\%$ ✓

$\{1, 2\} \Rightarrow \{3\}$ $2/4 = 50\%$

$\{1, 2\} \Rightarrow \{5\}$ $2/4 = 50\%$

$\{1, 3\} \Rightarrow \{2\}$ $2/4 = 50\%$

$\{1, 5\} \Rightarrow \{2\}$ $2/2 = 100\%$ ✓

$\{2, 3\} \Rightarrow \{1\}$ $2/4 = 50\%$

$\{2, 5\} \Rightarrow \{1\}$ $2/2 = 100\%$ ✓

没有包含信息的70%的

最终产生的规则 $\{5\} \Rightarrow \{1, 2\}$, $\{1, 5\} \Rightarrow \{2\}$, $\{2, 5\} \Rightarrow \{1\}$

算置信度时分子就是该子-项集的支持度, 分母是蕴涵式右例集合的支持度

