

## Assignment 1

Oct 12, 2015

Exercises 3, 8, 9, 11, 14, 16 in Chapter 6 of the textbook “Introduction to Data Mining”.

Students who don't have the textbook can download the pdf file named “ch6.pdf” from the QQ group for this course (Group#: 324742036).

Please complete the assignment using the following MS Word template and print on A4 papers. Please make sure your name and ID are correctly filled in your assignment.

Assignment due date: Oct 26, 2015, on class in classroom A12.

## 数据挖掘作业 1

## 频繁模式与关联分析

15S103163

宋博宇

$$3. (a). C(\phi \rightarrow A) = \frac{\sigma(\phi \cup A)}{\sigma(\phi)} = \frac{\sigma(A)}{|T|} = S(\phi \rightarrow A)$$

$$C(A \rightarrow \phi) = \frac{\sigma(\phi \cup A)}{\sigma(A)} = \frac{\sigma(A)}{\sigma(A)} = 1$$

$$(b). C_1 = \frac{\sigma(P \cup Q)}{\sigma(P)}$$

因为  $\sigma(P) \geq \sigma(P \cup Q) \geq \sigma(P \cup Q \cup R)$ 

$$C_2 = \frac{\sigma(P \cup Q \cup R)}{\sigma(P)}$$

所以  $C_1 \geq C_2$  $C_3 \geq C_2$ 

$$C_3 = \frac{\sigma(P \cup Q \cup R)}{\sigma(P \cup R)}$$

 $C_2$  是最小的置信度.

$$(c). S_1 = \frac{\sigma(P \cup Q)}{|T|}$$

因为  $S_1 = S_2 = S_3$ 

$$S_2 = \frac{\sigma(P \cup Q \cup R)}{|T|}$$

所以  $\sigma(P \cup Q) = \sigma(P \cup Q \cup R)$ 所以  $C_1 = C_2$ 

$$S_3 = \frac{\sigma(P \cup Q \cup R)}{|T|}$$

 $C_3 \geq C_1 = C_2$  $C_3$  是最大置信度, 或者三个都相同.

(d). 在这样的一个表格中:

A B

A

A B C

B C

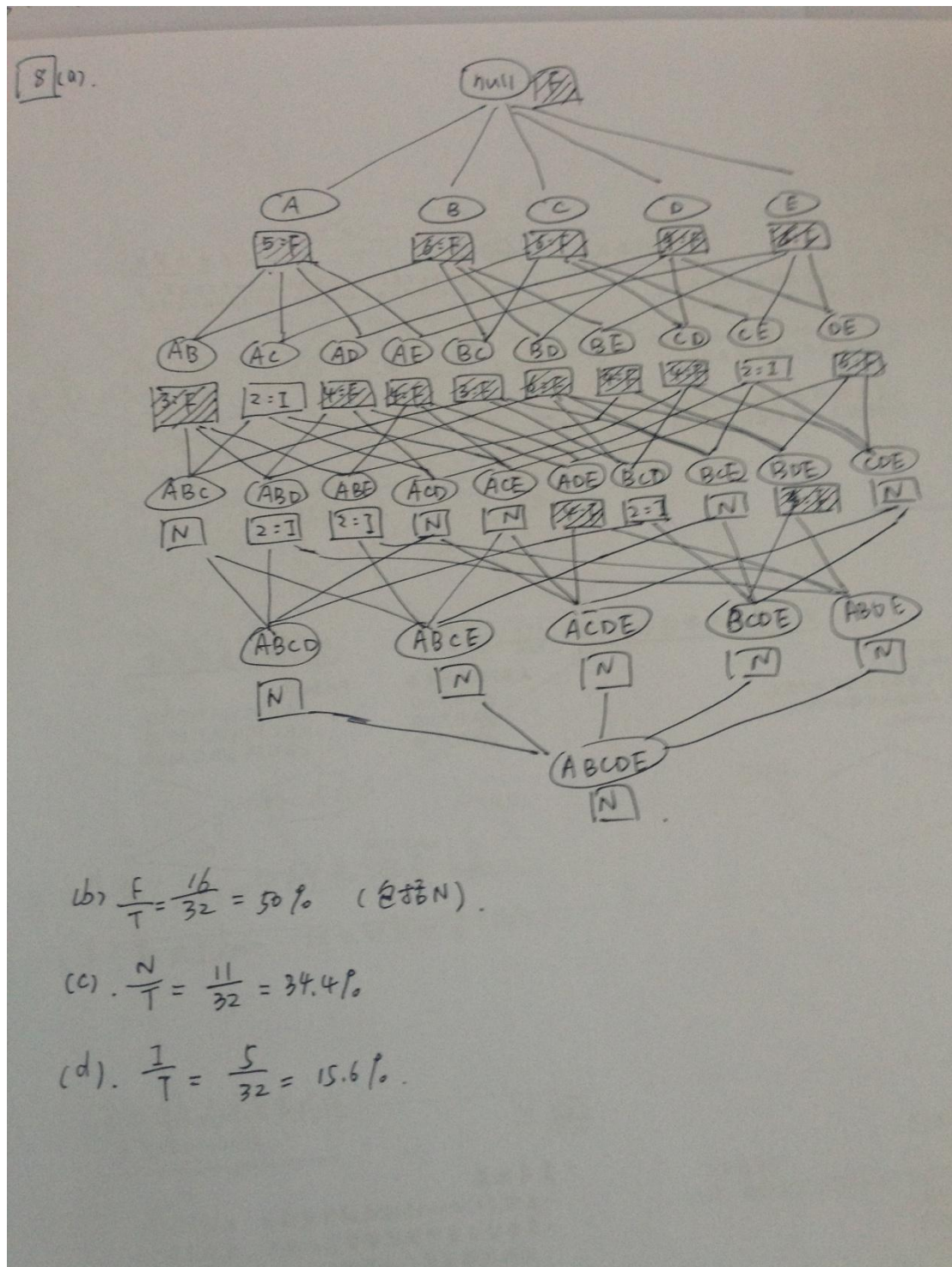
设  $\text{minconf} = \frac{1}{2}$ 

$$C(A \rightarrow B) = \frac{\sigma(A, B)}{\sigma(A)} = \frac{2}{3} > \frac{1}{2}$$

$$C(B \rightarrow C) = \frac{\sigma(B, C)}{\sigma(B)} = \frac{2}{3} > \frac{1}{2}$$

$$C(A \rightarrow C) = \frac{\sigma(A, C)}{\sigma(A)} = \frac{1}{3} < \frac{1}{2}$$

所以  $C(A \rightarrow C)$  有可能小于  $\text{minconf}$ .

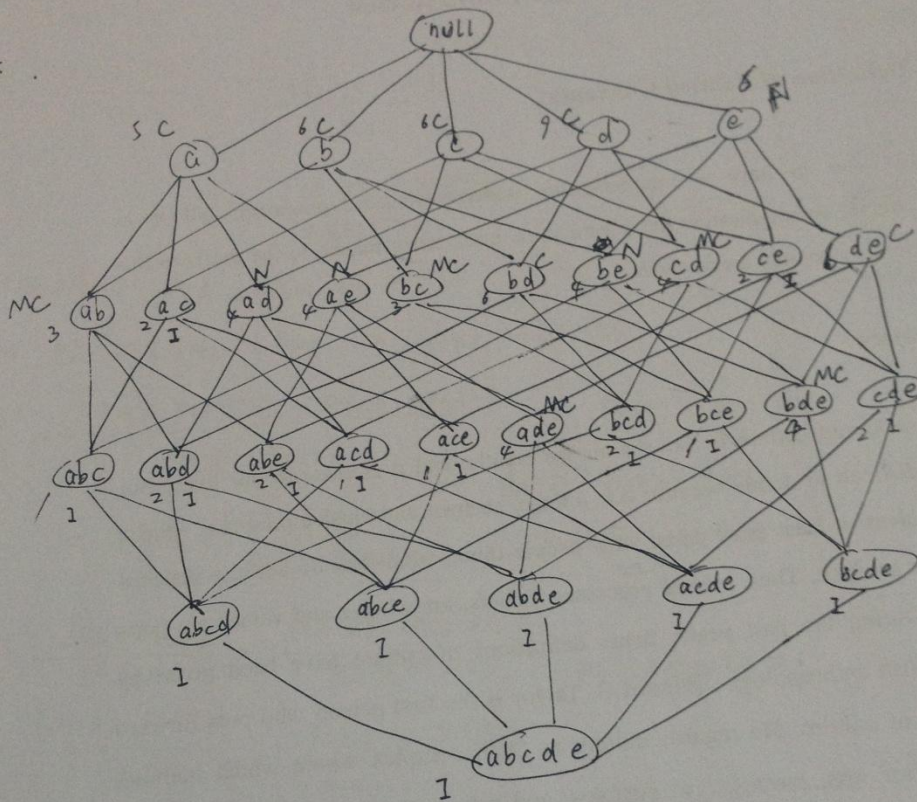




9. (a).  $L_1, L_3, L_5, L_9, L_{11}$ .

(b).  $\{145\}$   $\{158\}$   $\{458\}$   
 $L_1$   $L_3$   $L_3$

11.



14.

(a). e. 有频繁项集最长, 所以会产生最多频繁项集

(b). d. 没有项集个数大于 1000. 所以频繁项集个数为 0.

(c). e

(d). b. 最右侧长集表示有最大支持度.

(e). e. 支持度变化最广.

(a). 范围在  $[0, 1]$  之间.

$$M = \frac{P(B|A) - P(B)}{1 - P(B)}$$

$$\cancel{P(B|A) \geq P(B)} \quad P(A, B) \geq P(A)P(B)$$

最大值: 当  $P(B|A) = 1$  时  $M = \frac{1 - P(B)}{1 - P(B)} = 1$  即  $P(A, B) = P(A)$

最小值: 当  $P(B|A) = P(B)$  时  $M = 0$  即  $P(A, B) = P(A)P(B)$

(b).  $P(B|A) = \frac{P(A, B)}{P(A)}$

所以  $M = \frac{P(B|A) - P(B)}{1 - P(B)} = \frac{\frac{P(A, B)}{P(A)} - P(B)}{1 - P(B)}$

$$= \frac{P(A, B) - P(A)P(B)}{P(A)(1 - P(B))}$$

当  $P(A), P(B)$  不变时,  $P(A, B)$  增大则  $M$  增大.

(c). 由 (b) 可知  $P(A, B)$  和  $P(B)$  不变时,  $P(A)$  增大则  $M$  减小.

(d). 由 (b) 可知  $M = \frac{P(A) - P(A)P(B) + P(A, B) - P(A)}{P(A)(1 - P(B))}$   $P(A) \geq P(A, B)$

$$= 1 - \frac{P(A) - P(A, B)}{P(A)(1 - P(B))}$$

$P(A)$  和  $P(A, B)$  不变时  
 $P(B)$  增大  $M$  减小.

(e).  $M(A \rightarrow B) = \frac{P(B|A) - P(B)}{1 - P(B)} = \frac{P(A, B) - P(A)P(B)}{P(A)(1 - P(B))}$

所以不是对称的.

$$M(B \rightarrow A) = \frac{P(A|B) - P(A)}{1 - P(A)} = \frac{P(A, B) - P(A)P(B)}{P(B)(1 - P(A))}$$

(f). 0  $P(A, B) = P(A)P(B)$   $\leftarrow$  不依赖.

(g).  $M = \frac{P(A, B) - P(A)P(B)}{P(A) \cdot P(B)}$

所以  $S$  变化时

即  $\bar{B}$  变化,  $M$  也会变化

所以不是 null-invariant

|           | B | $\bar{B}$ |
|-----------|---|-----------|
| A         | P | q         |
| $\bar{A}$ | r | s         |

(h). 当  $A = \bar{A}$  时  $M = \frac{P(A, B) - P(A)P(B)}{P(A) \cdot P(B)}$

$$M = \frac{P - (P+q)(P+r)}{(P+q)(q+s)} \quad \text{当 } A = \bar{A} \text{ 时} \quad M = \frac{2P - 2(P+q)(2P+r)}{2(P+q)(2q+s)}$$

变化  
所以不是

(i) 反转时

|           | B | $\bar{B}$ |
|-----------|---|-----------|
| A         | s | q         |
| $\bar{A}$ | r | p         |

$$M = \frac{(s+r) - (s+q)(s+r)}{(s+q)(q+p)}$$

变化  
所以不对称