

19 Q4.

a	b	c	d	e	f
1		1			
	1	1	1		
1	1		1	1	
		1	1		
1	1	1	1		
		1	1	1	
1		1	1	1	
	1	1	1	1	

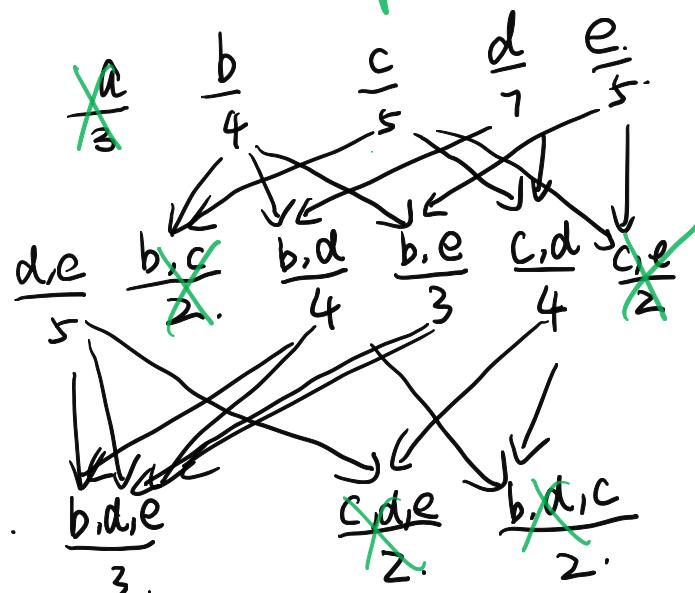
T1> maximal itemset.

直接超集都不频繁

T1> closed

~

~ 都没有与之相同的
sup.



a) T1> Find. sup≥3

maximal itemsets:

$$\{b, d, e\}, \{c, d\}$$

closed itemsets:

$$\{c\}, \{d\}, \{d, e\}, \{b, d\}, \{c, d\}, \{b, d, e\}$$

~~b, c, d, e~~

i) $b \rightarrow de$

$$\text{support: } \frac{s(d \cup e \cup b)}{N} = \frac{3}{8}$$

$$\text{confidence: } \frac{s(d \cup e \cup b)}{s(b)} = \frac{3}{4}$$

$$\text{lift: } \frac{CC(b \rightarrow de)}{s(de)} = \frac{3/4}{5} = \frac{3}{20}$$

b) i) Step 1:

$$\text{Set 1: } \{C_1 = 10 \mid 10\}.$$

$$\text{Set 2: } \{C_2 = 20 \mid 20\}.$$

$$\text{Set 3: } \{C_3 = 30 \mid 30, 40, 50, 60\} \Rightarrow C_3 = 45.$$

Step 2:

$$\text{Set 1: } \{C_1 = 10 \mid 10\}$$

$$\text{Set 2: } \{C_2 = 25 \mid 20, 30\} \Rightarrow C_2 = 25$$

$$\text{Set 3: } \{C_3 = 50 \mid 40, 50, 60\} \Rightarrow C_3 = 50$$

Step 3:

$$\text{Set 1: } \{C_1 = 10 \mid 10\}$$

$$\text{Set 2: } \{C_2 = 25 \mid 20, 30\}$$

$$\text{Set 3: } \{C_3 = 50 \mid 40, 50, 60\}$$

ii). $SSE = \sum \sum \text{dist}(c_i, x)^2$

$$= 0 + 5^2 + 5^2 + 10^2 + 10^2$$

$$= 250$$

$$\begin{aligned} BSS/SSB &= \sum |C_i| (C_i - \bar{C})^2 \\ &= 1 \times (10 - 35)^2 + 2 \times (25 - 35)^2 + 3 \times (50 - 35)^2 \\ &= 1500 \end{aligned}$$

$$TSS = \sum (x - \bar{C})^2$$

$$= 25^2 \times 2 + 15^2 \times 2 + 5^2 \times 2$$

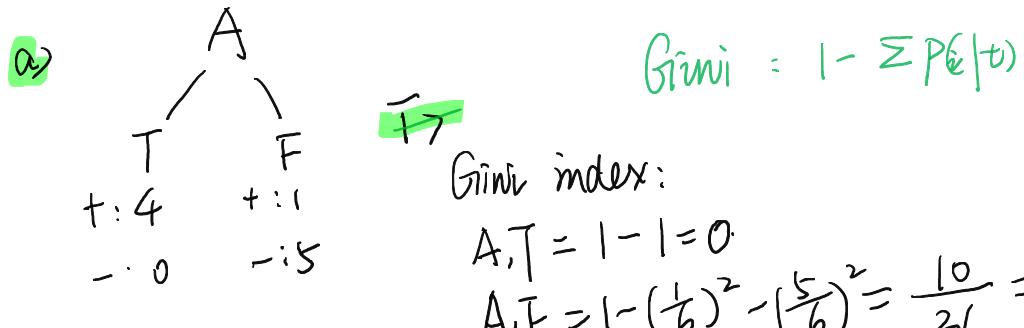
$$= 625 \times 2 + 225 \times 2 + 25$$

$$= 1250 + 450 + 50$$

$$= 1750$$

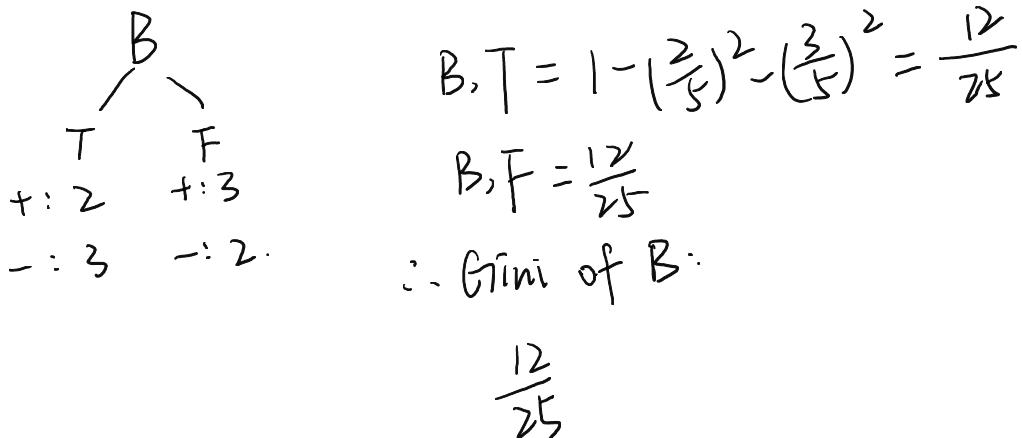
$$= SSE/WSS + SSB/BSS$$

18 Q4.



\therefore Gini of A:

$$\frac{6}{10} \times \frac{5}{18} = \frac{1}{6}$$



choose A.

i)

A:

	+	-
T	4	0
F	1	5

iii) Gain in Gini index

Overall Gini before splitting:

$$1 - (0.5)^2 - (0.5)^2 = 0.5$$

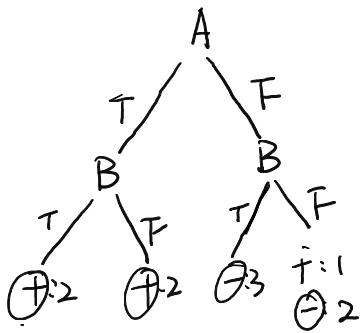
$$G_A = 0.5 - \frac{1}{6} = \frac{1}{3}$$

$$G_B = 0.5 - \frac{12}{25} = \frac{1}{50}$$

B:

	+	-
T	2	3
F	3	2

b>



C> ~~(add)~~

		+	-
		4 tp	0 fp
predict	+	1 fn	5 tn
	-		

$$acc = 0.9.$$

$$\text{precision} = \frac{\text{true positive}}{\text{tp} + \text{fp}}$$

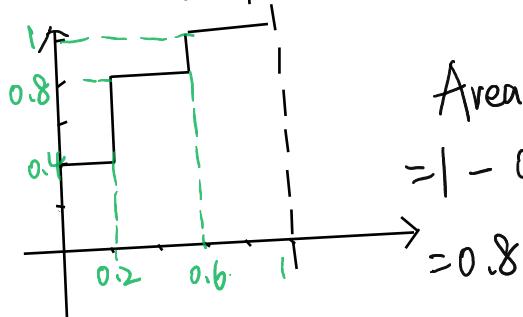
$$\text{recall} = \frac{0.8}{(\text{tp} + \text{fn})}$$

$$F_1 = \frac{2 \times 1 \times 0.8}{1+0.8}$$

$$= \frac{1.6}{1.8} = \frac{8}{9}$$

d>

	-	-	+	-	-	+	+	-	+	+
ID	X10	9	8	7	6	5	4	3	2	1
P		0.35	0.4	0.45						
TP	5	5	5	4	4	4	3	2	2	1
FP	5	4	3	3	2	1	1	1	0	0
TN	1	2	2	3	4	4	4	5	5	5
FN				1	1	1	2	3	3	4
TPR	1	1	1	0.8	0.8	0.8	0.6	0.4	0.4	0.2
FPR	1	0.8	0.6	0.6	0.4	0.2	0.2	0.2	0	0



Area under ROC

$$= 1 - 0.6 \times 0.2 - 0.4 \times 0.2$$

$$= 0.8$$

17 Q4

Step 1:

(a) $C_1 : 1100$
 $S_1 : \{0, 200, 300, 900, 1100\}$

$$C_1 \rightarrow 500.$$

$$C_2 : 1600$$

$$S_2 : \{1600\}$$

Step 2:

$$\frac{2100}{1050}$$

$$C_1 : 500$$

$$S_1 : \{0, 200, 300, 900\}$$

$$C_1 \rightarrow 350$$

$$C_2 : 1600$$

$$S_2 : \{1100, 1600\}$$

$$\frac{1400}{4} = 35$$

$$850.$$

$$C_2 \rightarrow 1350$$

Step 3

$$C_1 : 350$$

$$S_1 : \{0, 200, 300\}$$

$$C_1 = \frac{500}{3}$$

$$C_2 : 1350$$

$$S_2 : \{900, 1100, 1600\}$$

$$C_2 : 1200$$

Step 4.

Same as Step 3

(b) ✓

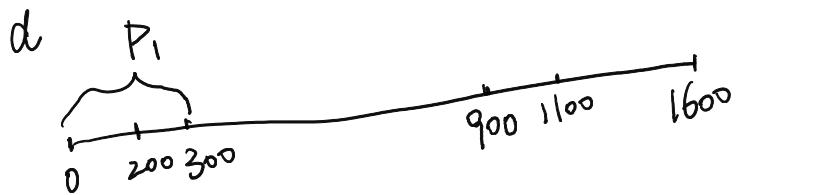
	P1	P2	P3	P4	P5	P6
P1	0	200	300	900	1100	1600
P2	200	0	100	700	9	14
P3	3	1	0	6	8	13
P4	9	7	6	0	2	7
P5	11	9	8	2	0	5
P6	16	14	13	7	5	0

	a _{ii}	b _{ii}	min _j {a _{ij} i ≠ j}	到 cluster 中其它 点的平均 dist
1	2.5	12		
2	1.5	10	0.85	
3	2	9		
4	4.5	22/3		
5	3.5	28/3		
6	6	43/3		

	a _{ii}	b _{ii}	min _j {a _{ij} i ≠ j}	到 cluster 中其它 点的平均 dist
1	2.5	12		
2	1.5	10	0.85	
3	2	9		
4	4.5	22/3		
5	3.5	28/3		
6	6	43/3		

$$b_{ii} - a_{ii} / \max(a_{ii}, b_{ii})$$

$$\frac{28}{3} - 3.5 = 9\frac{1}{3} - 3.5 = 4.5\frac{1}{3} = \frac{4.5 \times 3 + 1}{28} = \frac{29}{56}$$



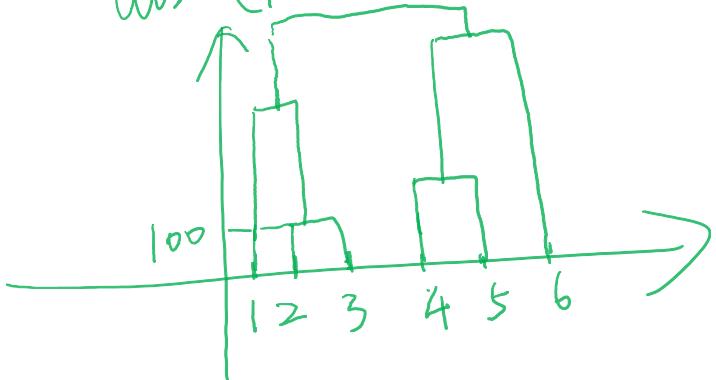
$$\textcircled{1} \quad \text{dist}(2, 3) = 100$$

$$\textcircled{2} \quad \text{dist}(4, 5) = 200$$

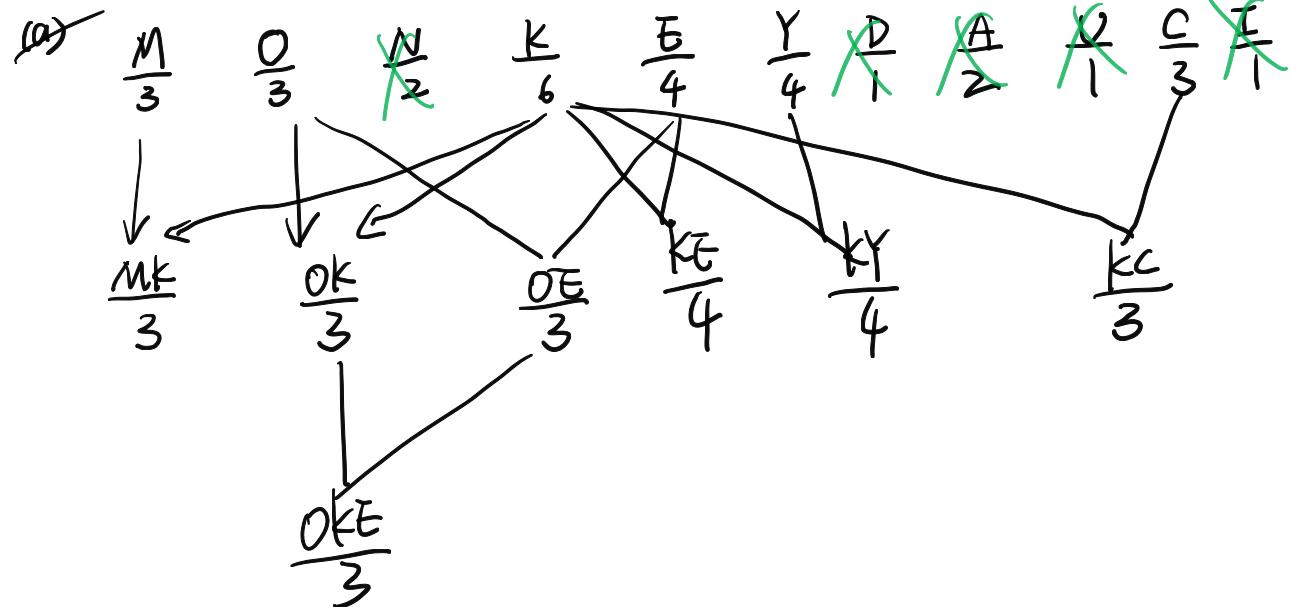
$$\textcircled{3} \quad \text{dist}(\{2, 3\}, 1) = 300 \quad \checkmark \rightarrow \{2, 3, 1\}$$

$$\text{dist}(\{4, 5\}, 6) = 700$$

$$\text{dist}(\{2, 3\}, \{4, 5\}) = 100 \times$$



16.
Q5



(b) $\text{con} = \frac{\sup\{E, C\}}{\sup\{E\}} = \frac{1}{4}.$

$\text{lift} = \frac{C\{E \rightarrow C\}}{S(C)} = \frac{1}{12}$

weakly related.

(c). $h(\{E, C\}) = \frac{\sup\{E, C\}}{\max(S(E), S(C))} = \frac{1}{\max(4, 3)} = \frac{1}{4}.$

不含。

$X \subset Y$

$\sup(X) \geq \sup(Y)$

$\max(S(X_1), S(X_2)) \leq \max(S(Y_1), S(Y_2), \dots)$

$\therefore h(X) \geq h(Y)$

\therefore it's anti-monotone.

$$\zeta(\{A, B\}) = \min [C(A \rightarrow B), C(B \rightarrow A)]$$

$$= \min \left[\frac{\sup(A, B)}{\sup(A)}, \frac{\sup(A, B)}{\sup(B)} \right]$$

$$= \frac{\sup(A, B)}{\max(\sup(A), \sup(B))}$$

$$C(\{A, B, C\}) = \min [C(A \rightarrow BC), C(B \rightarrow AC), C(C \rightarrow AB)]$$

$$= \frac{\sup(A, B, C)}{\max(\sup(A), \sup(B), \sup(C))}$$

$$\{A, B\} \subset \{A, B, C\}$$

$$\therefore \sup(A, B) \geq \sup(A, B, C)$$

$$\max(\sup(A), \sup(B)) \leq \max(\sup(A), \sup(B), \sup(C))$$

$$\therefore \overline{\cap}(A, B) \geq \overline{\cap}(A, B, C)$$

\therefore anti-monotone

15 Ex. 7 on Page 37
~~as~~

