

7.5.

Step one: For each frequent item, find its conditional pattern base and then its conditional FP-tree.

Repeat step one until resulting FP tree is empty or it only has one path. Then we store both frequent and infrequent patterns.

Step two: for each frequent itemset X , we scan infrequent item sets Z , which contain X .

Let $W = Z \setminus X$, i.e.: W contain elements in Z but not in X . If W is also a frequent itemset.

Step three: calculate $(P(X|W) + P(W|X))/2$

~~# ①~~ whether W and X must be negatively correlated depend on $(P(X|W) + P(W|X))/2$

$$= \left(\frac{\text{sup}(Z)}{\text{sup}(W)} + \frac{\text{sup}(Z)}{\text{sup}(X)} \right) / 2. \text{ note that}$$

$\text{sup}(X)$, $\text{sup}(Y)$ and $\text{sup}(Z)$ are known aprioris steps.

7.9. The distance measure Pat-Dist is a valid distance metric. It has following properties:

(a). $\text{Pat-Dist}(P_1, P_2) > 0, \forall P_1 \neq P_2$

(b). $\text{Pat-Dist}(P_1, P_2) = 0, \forall P_1 = P_2$

(c). $\text{Pat-Dist}(P_1, P_2) = \text{Pat-Dist}(P_2, P_1)$

(d). $\text{Pat-Dist}(P_1, P_2) + \text{Pat-Dist}(P_2, P_3) \geq \text{Pat-Dist}(P_1, P_3)$

Proof of (a): $\forall P_1 \neq P_2$, since P_1, P_2 are chosen patterns,

$$0 \leq \frac{|T(P_1) \cap T(P_2)|}{|T(P_1) \cup T(P_2)|} < 1 \Rightarrow \text{Pat-Dist}(P_1, P_2) > 0.$$

Proof of (b): $\forall P_1 = P_2 \quad \therefore |T(P_1) \cap T(P_2)| = |T(P_1) \cup T(P_2)|$
 $= |T(P_1)| = |T(P_2)| \quad \therefore \text{Pat-Dist}(P_1, P_2) = 0$.

(c). is obvious, since the formula is symmetric.

(d). Since $(T(P_1) \cap T(P_2)) \cup (T(P_1) \cap T(P_3)) \subseteq T(P_1)$

we have

$$\therefore A \cup B = A + B - A \cap B$$

$$\therefore |T(P_1) \cap T(P_2)| + |T(P_1) \cap T(P_3)| - |T(P_1) \cap T(P_2) \cap T(P_3)| \leq |T(P_1)| \Rightarrow b_1 + c_1 - d_1 \leq a \quad (*)$$

$$(d) \Leftrightarrow 1 - \frac{b_1}{a+b_2} + 1 - \frac{c_1}{a+c_2} \geq 1 - \frac{d_1+d_2}{b_1+b_2+c_1+c_2-d_1-d_2}$$

$$\Leftrightarrow \frac{b_1}{a+b_2} + \frac{c_1}{a+c_2} \leq 1 + \frac{d_1+d_2}{b_1+b_2+c_1+c_2-d_1-d_2}$$

$$\therefore 1 + \frac{d_1+d_2}{b_1+b_2+c_1+c_2-d_1-d_2} \geq 1 + \frac{d_1+d_2}{b_1+b_2+c_1+c_2-d_1}$$

$$(d_2 \geq 0) \geq 1 + \frac{d_1+d_2}{b_1+b_2+c_1+a} \quad (*) = \frac{a+b_2+c_1+d_1}{a+b_2+c_1}$$

$$\geq \frac{b_1+c_1+b_2+c_2}{a+b_2+c_1} \quad (*) = \frac{b_1+c_1}{a+b_2+c_1} + \frac{b_2+c_1}{a+b_2+c_1}$$

$$\geq \frac{b_1}{a+b_2} + \frac{c_1}{a+c_2} \quad (a+b_2 \geq b_1, c_1 \geq 0) \quad (a+c_2 \geq c_1, c_1 \geq 0).$$

Hence the (d) holds.

7.10: In order to compress a large set of patterns to a small representatives, we can apply clustering method. We use the concept of S -cluster. Given a transaction database, a minimum support min-sup, the cluster quality measure S , the pattern compression problem is to find a set of representative patterns R such that for each frequent pattern $P \in R$ (w.r.t min-sup), there's a representative pattern $P_r \in R$, which covers P , and $|R|$ is minimized.

8.3: If we directly prune decision tree and convert to rules, we would remove the subtree entirely. But if we prune after converting decision tree to rules, we may remove any precondition of it. The latter results in loss information loss.

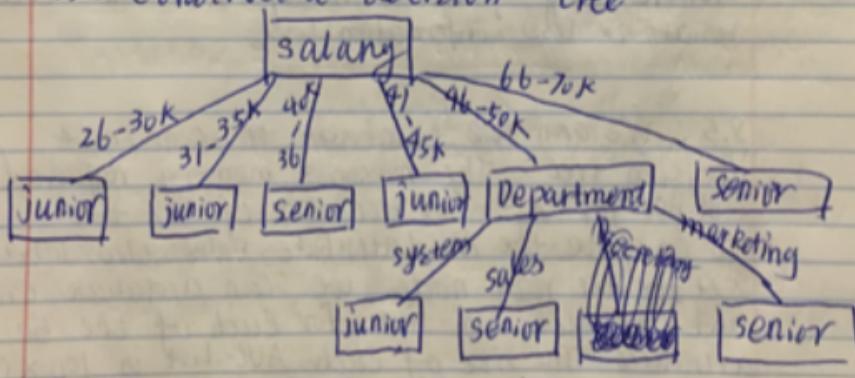
8.5: We can use RainForest to construct decision tree. The maximal memory required will be for AVC-set for root of the tree. To compute the AVC(attribute, value, class-label) set for the root node, we scan database once and construct AVC-list for each of the 50 attributes. The size of each AVC-list is $100 \times C$, C is the number of class labels. Then size of AVC-set is $50 \times 100 \times C$, which can be easily fit into 512 MB memory given a reasonable C .

For other nodes, AVC-sets are done in similar ways but they will be smaller since they will consider less number of attributes. To reduce number of scan for the original data set, we can parallelly compute AVC-set for nodes at same level in the tree.

8.7.

(a). Basic decision tree algorithm must be modified to consider the count of generalized data tuple. ①. Count is considered to calculate the most common class among all the tuples. ②. the count of each tuple must be used when calculating attribute selection measure, such as information gain.

(b). Construct a decision tree



(c). According to Bayesian Theorem.

Given such an observation data X , posterior probability of a Hypothesis H

actual	Predict		P
	TP	FN	
	FP	TN	N
	P'	N'	

$$P(H|X) = \frac{P(X|H) P(H)}{P(X)}$$

let H: be X is senior

$$P(\text{Senior}|X) = \frac{P(X|\text{senior}) P(\text{senior})}{P(X)} = 0$$

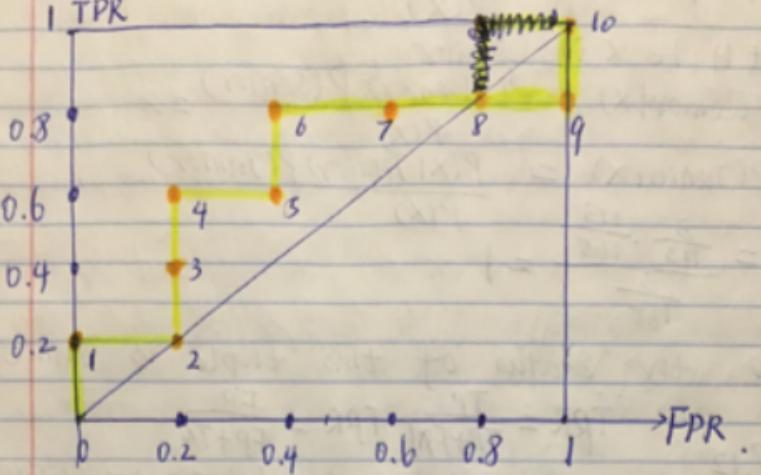
$$P(\text{Junior}|X) = \frac{P(X|\text{Junior}) P(\text{Junior})}{P(X)} \\ = \frac{\frac{3}{113} \cdot \frac{113}{165}}{\frac{3}{165}} = 1$$

So, the status of this tuple is Junior.

$$8.12 \quad TPR = \frac{TP}{TP+FN} \quad FPR = \frac{FP}{FP+TN}$$

Tuple #	Class	prob	TP	FN	FP	TN	TPR / FPR
1	P	0.95	1	4	0	5	0.2 0
2	n	0.85	1	4	1	4	0.2 0.2
3	P	0.78	2	3	1	4	0.4 0.2
4	P	0.66	3	2	1	4	0.6 0.2
5	n	0.60	3	2	2	3	0.6 0.4
6	P	0.55	4	1	2	3	0.8 0.4
7	n	0.53	4	1	3	2	0.8 0.6
8	n	0.52	4	1	4	1	0.8 0.8
9	n	0.51	4	1	5	0	0.8 1.0
10	P	0.4	5	0	5	0	1.0 1.0

ROC Curve:



8.14 :

there's no significant

Let H_0 be Null hypothesis : $\overline{err}(M_1) = \overline{err}(M_2)$

H_1 M_1 is better than M_2 . $\overline{err}(M_1) < \overline{err}(M_2)$

$$t\text{-statistic} = \frac{\overline{err}(M_1) - \overline{err}(M_2)}{\sqrt{\text{Var}(\overline{err}(M_1) - \overline{err}(M_2)) / K}}$$

$$\text{where } \text{Var}(\overline{err}(M_1) - \overline{err}(M_2)) = \frac{1}{K} \sum_{i=1}^K [(\overline{err}(M_1)_i - \overline{err}(M_2)_i - (\overline{err}(M_1) - \overline{err}(M_2))]^2$$

$$\therefore \text{Var}(\overline{err}(M_1) - \overline{err}(M_2)) = \frac{1}{10} (1.65^2 + 11.25^2 + (-8.15)^2 + (-5.45)^2 + 3.85^2 + 14.15^2 + (-0.85)^2 + 0.15^2 + (-1.15)^2 + (-15.45)^2) \\ = 68.1225$$

$$t = \frac{27.72 - 21.27}{\sqrt{68.1225 / 10}} = 2.47$$

For t -distribution with $10 - 1 = 9$ degree of freedom

(two tail)

K for 1% significance level is 3.25
since $2.47 < 3.25$, i.e. $2.47 < 3.25$

we fail to reject null hypothesis, hence,
there's no significant better model
among / between M_1 and M_2 .