

Exercise 11.1 Tan, Chapter 3

Ex 2]

2- (1) Gini index for the overall collection of training examples.

$$\Rightarrow \text{Gini index} = 1 - \sum p_i t_i^2$$

$$= 1 - \left[\left(\frac{10}{20} \right)^2 + \left(\frac{10}{20} \right)^2 \right] = 1 - \left[\frac{1}{4} + \frac{1}{4} \right] = 1 - \frac{2}{4} = \underline{\underline{0.5}}$$

2- (2) Compute the Gini index for the customer ID attribute

$$\Rightarrow 1 - \left[\left(\frac{0}{1} \right)^2 + \left(\frac{1}{1} \right)^2 \right] \Rightarrow 1 - 1 \Rightarrow \underline{\underline{0}}$$

2- (3) Compute the Gini index for the gender attribute

for females,

$$\text{Gini index} = 1 - \left[\left(\frac{6}{10} \right)^2 + \left(\frac{4}{10} \right)^2 \right]$$

$$= 1 - \cancel{\left(\left(0.3 \right)^2 + \left(0.2 \right)^2 \right)}$$

$$\Rightarrow 1 - \left(0.6 \right)^2 + \left(0.4 \right)^2$$

$$\Rightarrow 1 - (0.36 + 0.16)$$

$$\Rightarrow \underline{\underline{0.48}}$$

for males,

$$\begin{aligned}\text{Gini index} &= 1 - \left(\left(\frac{4}{10}\right)^2 + \left(\frac{6}{10}\right)^2 \right) \\ &= 1 - 0.52 \\ &= 0.48.\end{aligned}$$

Gini index for gender would be,

$$\begin{aligned}\left(\frac{10}{20}\right) * 0.48 + \left(\frac{10}{20}\right) * 0.48 \\ \Rightarrow \underline{\underline{0.48}}\end{aligned}$$

2-(4) Compute the Gini index for the car type attribute using multiway split

\Rightarrow Family,

$$\begin{aligned}\text{Gini Index} &= 1 - \left[\left(\frac{1}{4}\right)^2 + \left(\frac{3}{4}\right)^2 \right] \\ &\Rightarrow 1 - \left((0.25)^2 + (0.75)^2 \right) \\ &\Rightarrow 0.375\end{aligned}$$

Luxury,

$$\begin{aligned}G.I &= 1 - \left(\left(\frac{1}{8}\right)^2 + \left(\frac{7}{8}\right)^2 \right) \\ &\Rightarrow 0.219\end{aligned}$$

Spots,

$$G.I = 1 - \left(\left(\frac{8}{8}\right)^2 + \left(\frac{0}{8}\right)^2 \right)$$

$$\Rightarrow 1 - 1 = \underline{\underline{0}}$$

Q Gini index of car attribute is,

$$\left(\frac{4}{20}\right) * 0.375 + \left(\frac{8}{20}\right) * 0.219 + \frac{8}{20} * 0$$

$$\Rightarrow \underline{\underline{0.163}}$$

2-(S) Compute Gini Index of shirt size attribute using multiway split.

\Rightarrow Small,

$$G.I = 1 - \left(\left(\frac{3}{5}\right)^2 + \left(\frac{2}{5}\right)^2 \right)$$

$$\Rightarrow 1 - \left((0.6)^2 + (0.4)^2 \right)$$

$$\Rightarrow \underline{\underline{0.48}}$$

Medium,

$$G.I = 1 - \left(\left(\frac{3}{7}\right)^2 + \left(\frac{4}{7}\right)^2 \right)$$

$$\Rightarrow 1 - \left((0.43)^2 + (0.57)^2 \right)$$

$$\Rightarrow \underline{\underline{0.49}}$$

Large,

$$GI = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right]$$

$$\Rightarrow 1 - \left((0.5)^2 + (0.5)^2 \right)$$

$$= \underline{\underline{0.5}}.$$

Extra Large,

$$GI = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right]$$

$$= 1 - \left((0.5)^2 + (0.5)^2 \right)$$

$$= 0.5.$$

Gini index for shirt attribute is given by,

$$\left(\frac{5}{20} \right) * 0.48 + \left(\frac{7}{20} \right) * 0.49 + \left(\frac{4}{20} \right) * 0.5 + \left(\frac{4}{20} \right) * 0.5$$
$$\Rightarrow \underline{\underline{0.4915}}.$$

2-(b) which attribute is better, gender, car or shirt?

\Rightarrow Among the three attributes, car type is better as it has the lowest Gini index of all i.e., 0.163 whereas the other 2 have 0.48 (gender) and (~~0.49~~) 0.49 (shirt size).

2-(7) Why customer ID should not be used as the attribute test condition even though it has the lowest Gini.

⇒ Because each attribute has unique customer id. And hence the attribute has no predictive power.

Ex 3

3-(1) Entropy =,

$$E(S) = \sum_{i=1}^c -P_i \log_2 P_i$$

$$= \left[\left(\frac{4}{9} \right) \log_2 \left(\frac{4}{9} \right) + \left(\frac{5}{9} \right) \log_2 \left(\frac{5}{9} \right) \right]$$

$$= -(-0.519) + (-0.471)$$

$$= \underline{\underline{0.991}}$$

3-(2) For a_1 , $T = (+ve \rightarrow 3, -ve \rightarrow 1) \mid F = (+ve \rightarrow 1, -ve \rightarrow 4)$

For a_2 , $T = (+ve \rightarrow 2, -ve \rightarrow 3) \mid F = (+ve \rightarrow 2, -ve \rightarrow 2)$

Entropy for a_1

$$\Rightarrow \frac{4}{9} \left[\cancel{- \left[\left(\frac{3}{4} \right) \log_2 \left(\frac{3}{4} \right) + \left(\frac{1}{4} \right) \log_2 \left(\frac{1}{4} \right) \right]} + \cancel{\frac{5}{9} \left[- \left[\left(\frac{1}{5} \right) \log_2 \left(\frac{1}{5} \right) + \left(\frac{4}{5} \right) \log_2 \left(\frac{4}{5} \right) \right]} \right]$$

⇒ Entropy for $a_1(T)$

$$\Rightarrow - \left[\left(\frac{3}{4} \right) \log_2 \left(\frac{3}{4} \right) + \left(\frac{1}{4} \right) \log_2 \left(\frac{1}{4} \right) \right] \Rightarrow \underline{\underline{0.81128}}$$

Entropy for $a_1(F)$

$$E = - \left[\left(\frac{1}{5}\right) \log_2 \left(\frac{1}{5}\right) + \left(\frac{4}{5}\right) \log_2 \left(\frac{4}{5}\right) \right]$$

$$\Rightarrow [0.721]$$

Entropy for $a_2(T)$

$$E = - \left[\left(\frac{2}{5}\right) \log_2 \left(\frac{2}{5}\right) + \left(\frac{3}{5}\right) \log_2 \left(\frac{3}{5}\right) \right]$$

$$= [0.970]$$

Entropy for $a_3(F)$

$$E = - \left[\left(\frac{2}{5}\right) \log_2 \left(\frac{2}{5}\right) + \left(\frac{2}{5}\right) \log_2 \left(\frac{2}{4}\right) \right]$$

$$\Rightarrow [1]$$

Information gain

$$\text{for } a_1 \Rightarrow (0.99) - \left[\left(\frac{4}{9}\right) * (0.81128) + \left(\frac{5}{9}\right) * (0.721) \right]$$

$$\Rightarrow 0.229$$

$$\text{for } a_2 \Rightarrow (0.99) - \left[\left(\frac{5}{9}\right) * (0.970) + \left(\frac{4}{9}\right) * 1 \right]$$

$$\underline{\underline{= 0.007}}$$

a_3	class label
1.0	+
6.0	+
5.0	-
4.0	+
7.0	-
3.0	-
8.0	-
7.0	+
5.0	-

Sorted a_3	class label	Split point	Entropy	Info gain
1.0	+	2.0	0.848	0.142
3.0	-	3.5	0.988	0.0024
4.0	+	4.5	0.918	0.072
5.0	-	5.5	0.983	0.0072
5.0	-			
6.0	+	6.5	0.972	0.018
7.0	+			
7.0	-	7.5	0.888	0.102

~~The~~
The best split for a_3 occurs at split point equals to 2.

Split 1 split point 2.0.

$$\leq(E) = - \left[\left(\frac{1}{4}\right) \log_2 \left(\frac{1}{4}\right) + 0 \log_2 0 \right] \\ = 0.$$

$$>(E) = - \left[\left(\frac{3}{8}\right) \log_2 \left(\frac{3}{8}\right) + \left(\frac{5}{8}\right) \log_2 \left(\frac{5}{8}\right) \right] \\ \Rightarrow -((-0.530) + (-0.423)) \\ = 0.954.$$

Weighted average = $\left[\left(\frac{1}{9}\right) * 0 \right] + \left(\frac{8}{9}\right) * 0.954 \\ = 0.848$

Information gain $\Rightarrow 0.991 - 0.848 = 0.143$

Split 2 point 3.5

$$\leq(E) = - \left[\left(\frac{1}{2}\right) \log_2 \left(\frac{1}{2}\right) + \left(\frac{1}{2}\right) \log_2 \left(\frac{1}{2}\right) \right] \\ = -(-0.5) + (-0.5) = 1.$$

$$>(E) = - \left[\left(\frac{3}{7}\right) \log_2 \left(\frac{3}{7}\right) + \left(\frac{4}{7}\right) \log_2 \left(\frac{4}{7}\right) \right] \\ = -(-0.523) + (-0.461) \\ = 0.985$$

Weighted avg = $\left[\left(\frac{2}{9}\right) * 1 + \left(\frac{7}{9}\right) * 0.985 \right] = 0.988.$
Information gain = 0.0024 //

Split point 4,5

$$\leq (E) - \left[\left(\frac{2}{3} \right) \log_2 \left(\frac{2}{3} \right) + \left(\frac{1}{3} \right) \log_2 \left(\frac{1}{3} \right) \right]$$

$$= 0.918$$

$$>(E) = - \left[\left(\frac{2}{6} \right) \log_2 \left(\frac{2}{6} \right) + \left(\frac{4}{6} \right) \log_2 \left(\frac{4}{6} \right) \right]$$

$$= 0.918.$$

$$\text{Weighted avg} = \left[\left(\frac{3}{9} \right) * (0.918) + \left(\frac{6}{9} \right) * (0.918) \right]$$

$$= 0.918.$$

$$\text{Information gain} \Rightarrow 0.991 - 0.918 = \underline{\underline{0.072}}$$

Split point 5,5

$$\leq (E) = - \left[\left(\frac{2}{5} \right) \log_2 \left(\frac{2}{5} \right) + \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right]$$

$$= 0.970$$

$$>(E) = - \left[\left(\frac{2}{4} \right) \log_2 \left(\frac{2}{4} \right) + \left(\frac{2}{4} \right) \log_2 \left(\frac{2}{4} \right) \right]$$

$$= 1.$$

$$\text{Weighted avg} = \left[\left(\frac{5}{9} \right) * 0.970 \right] + \left[\left(\frac{4}{9} \right) * 1 \right]$$
$$= 0.983$$

$$\text{Information gain} = 0.991 - 0.983 = \underline{\underline{0.0071}}$$

split point 6.5

$$\leq(E) = - \left[\left(\frac{3}{6} \right) \log_2 \left(\frac{3}{6} \right) + \left(\frac{3}{6} \right) \log_2 \left(\frac{3}{6} \right) \right] = 1.$$

$$>(E) = - \left[\left(\frac{1}{3} \right) \log_2 \left(\frac{1}{3} \right) + \left(\frac{2}{3} \right) \log_2 \left(\frac{2}{3} \right) \right] = 0.918.$$

$$\text{weighted avg} \Rightarrow \left[\left(\frac{6}{9} \right) * 1 \right] + \left[\left(\frac{3}{9} \right) * 0.918 \right] = 0.972.$$

$$\text{Information gain} = 0.991 - 0.972 = 0.019$$

split point 7.5

$$\leq(E) = - \left[\left(\frac{4}{8} \right) \log_2 \left(\frac{4}{8} \right) + \left(\frac{4}{8} \right) \log_2 \left(\frac{4}{8} \right) \right] = 1$$

$$>(E) = - \left[\left(\frac{0}{1} \right) \log_2 \left(\frac{0}{1} \right) + \left(\frac{1}{1} \right) \log_2 \left(\frac{1}{1} \right) \right] = 0.$$

$$\text{weighted avg} = \left[\left(\frac{8}{9} \right) * 1 \right] + \left[\left(\frac{1}{9} \right) * 0 \right] = 0.889$$

$$\text{Information gain} = 0.991 - 0.889 = 0.102.$$

3-(4) According to information gain,
 a_1 produces the best split as it has
higher gain

3-(5)

For $a_1(T)$

$$\Rightarrow 1 - \left[\frac{3}{4}, \frac{1}{4} \right] = 1 - \frac{3}{4} = 0.25$$

$a_1(F)$

$$\Rightarrow 1 - \left[\frac{1}{5}, \frac{4}{5} \right] = 1 - \frac{4}{5} = 0.2$$

Avg. error rate $\Rightarrow \frac{4}{9} * 0.25 + \frac{5}{9} * 0.2$
 $\Rightarrow \underline{\underline{0.222}}$

For $a_2(T)$

$$\Rightarrow 1 - \left[\frac{2}{5}, \frac{3}{5} \right] = 1 - \frac{3}{5} = 0.4$$

$a_2(F)$

$$\Rightarrow 1 - \left[\frac{2}{4}, \frac{2}{4} \right] = 1 - \frac{2}{4} = 0.5$$

Avg. error rate $= \frac{5}{9} * 0.4 + \frac{4}{9} * 0.5$
 $\Rightarrow 0.222 + 0.222$
 $\underline{\underline{= 0.444}}$

According to classification error rate, the best split is a₁ due to low classification error rate. The classification error depicts the accuracy of the sample set; the higher the classification error the more error the sample set contains.

$$3-(6) \\ \text{for } a_1(T), \\ \Rightarrow 1 -$$

$$\underline{3-(6)} \\ \text{for } a_1 \Rightarrow$$

$$\text{Gini Index} = \frac{4}{9} \left[1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] \right] + \frac{5}{9} \left[1 - \left[\left(\frac{1}{5} \right)^2 + \left(\frac{4}{5} \right)^2 \right] \right] \\ = 0.3444.$$

$$\text{for } a_2 \Rightarrow$$

$$\text{Gini Index} = \frac{4}{9} \left[1 - \left[\left(\frac{2}{8} \right)^2 + \left(\frac{3}{5} \right)^2 \right] \right] + \frac{4}{9} \left[1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] \right] \\ = 0.4889.$$

Since gini index of a₁ is small, it has the better split.

Ex 5

5 (a). The contingency tables after splitting on attributes A and B are:-

$A=T$	$A=F$	$B=T$	$B=F$
+	4 0	3 1	
-	3 3	1 5	

overall entropy before splitting ,

$$E_{\text{original}} = - (0.4 \log 0.4 + 0.6 \log 0.6)$$

$$= 0.97$$

The information gain after splitting on A.

$$E_A(T) = - \left(\left(\frac{4}{7}\right) \log_2 \left(\frac{4}{7}\right) + \left(\frac{3}{7}\right) \log_2 \left(\frac{3}{7}\right) \right)$$

$$= 0.985$$

$$E_A(F) = - \left(\left(\frac{3}{3}\right) \log_2 \left(\frac{3}{3}\right) + \left(\frac{0}{3}\right) \log_2 \left(\frac{0}{3}\right) \right)$$

$$= 0.$$

Info gain Δ , $\Rightarrow E_{\text{orig}} - (7/10) E_A(T) - (3/10) E_A(F)$

$$\Rightarrow 0.2813$$

The information gain after splitting on B is

$$E_B(T) = - \left(\left(\frac{3}{4}\right) \log_2 \left(\frac{3}{4}\right) + \left(\frac{1}{4}\right) \log_2 \left(\frac{1}{4}\right) \right) = 0.811$$

$$E_B(F) = - \left(\left(\frac{1}{6}\right) \log_2 \left(\frac{1}{6}\right) + \left(\frac{5}{6}\right) \log_2 \left(\frac{5}{6}\right) \right) = 0.650$$

$$\Delta = E_{\text{orig}} - (4/10) E_B(T) - (6/10) E_B(F) = 0.2565$$

Therefore, attribute A will be chosen to split the mode.

5(b)

The overall gini before splitting is

$$G_{\text{orig}} = 1 - 0.4^2 - 0.6^2 = 0.48$$

The gain in gini after splitting on A,

$$G_A(T) = 1 - \left(\frac{4}{7}\right)^2 - \left(\frac{3}{7}\right)^2 = 0.489$$

$$G_A(F) = 1 - \left(\frac{3}{3}\right)^2 - \left(\frac{0}{3}\right)^2 = 0.$$

$$\Delta = 0.48 - \left(\frac{7}{10}\right) * 0.489 - \left(\frac{3}{10}\right) * 0 \\ = 0.1371$$

The gain in gini after splitting on B,

$$G_B(T) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = 0.375$$

$$G_B(F) = 1 - \left(\frac{1}{6}\right)^2 - \left(\frac{5}{6}\right)^2 = 0.2778$$

$$\Delta = 0.48 - \left(\frac{4}{10}\right) * 0.375 - \left(\frac{6}{10}\right) * 0.2778 \\ = 0.1633$$

Therefore attribute B will be chose to split the mode.

5-(c) Yes, even though these measures have similar range and monotonous behavior, their respective gains, which are scaled differences of the measures, do not necessarily behave in the same way, as in results from part (a) & (b).

1.2 Tan, Chapter 4.

Question 18.

18-(1) Consider 100 total values, from this, 50 \rightarrow training set & the rest are used for testing

The error rate is the percent of incorrect classifications. Now for this case, for the records of 50% positive and 50% negative, the decision tree is predicting that all the records are positive.

So,

		Predicted	
		negative	positive
Negative(AV)	Negative	0	25
	Positive	0	25

$$\text{error rate} = \frac{(\text{False+ve} + \text{False-ve})}{(\text{T+ve}) + \text{F.P} + \text{T.N} + \text{F.N}} \Rightarrow \frac{25 + 0}{25 + 0 + 0 + 25}$$
$$\Rightarrow 0.5 = \underline{\underline{50\%}}$$

18-(2) positive class with probability 0.8
negative class with probability 0.2

		predicted	
		positive PV	Negative PV
Actual	positive AV	20	5
	negative AV	20	5

$$\text{Error rate} = (FP + FN) / (TP + TN + FP + FN)$$

$$\Rightarrow \frac{20 + 5}{20 + 5 + 20 + 5} \Rightarrow \frac{25}{50} = 0.5 \\ = 50\%$$

18-(3) from the set of 50 records,

(33.33) $\frac{2}{3} * 50 \rightarrow$ belongs to \rightarrow +ve class.

(16.66) $\frac{1}{3} * 50 \rightarrow$ belongs to \rightarrow -ve class

		Positive PV	Negative PV
Actual	positive AV	33.33	0
	Negative AV	16.66	0

$$\text{Error rate} \Rightarrow \frac{16.66 + 0}{33.33 + 0 + 16.66 + 0} \Rightarrow \frac{16.66}{49.99} \\ \Rightarrow 0.333 \Rightarrow \underline{\underline{33\%}}$$

18-(v) From 50 records,

Actual +ve \rightarrow 33.33

Actual -ve \rightarrow 16.66.

Predicted +ve \rightarrow 22.22.

Predicted -ve \rightarrow 5.55.

	positive PV	Negative PV
positive AV	22.22	11.11
Negative AV	11.11	5.55

$$\text{Error rate} = \frac{11.11 + 11.11}{22.22 + 5.55 + 11.11 + 11.11}$$

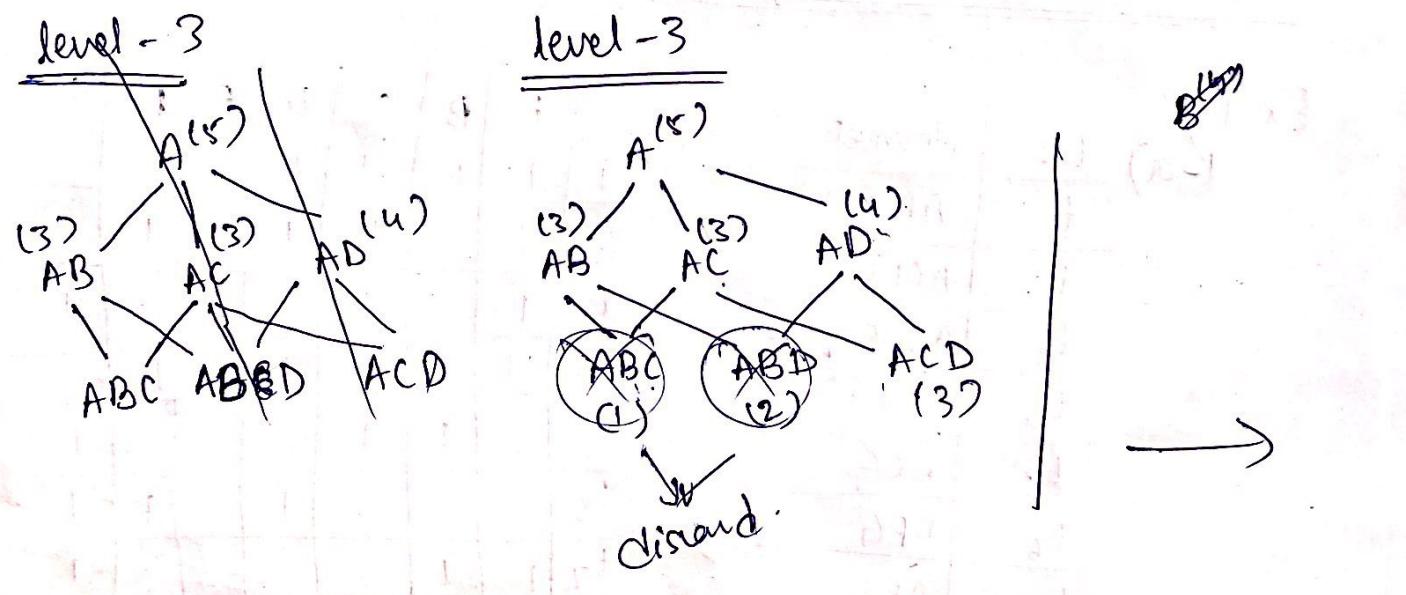
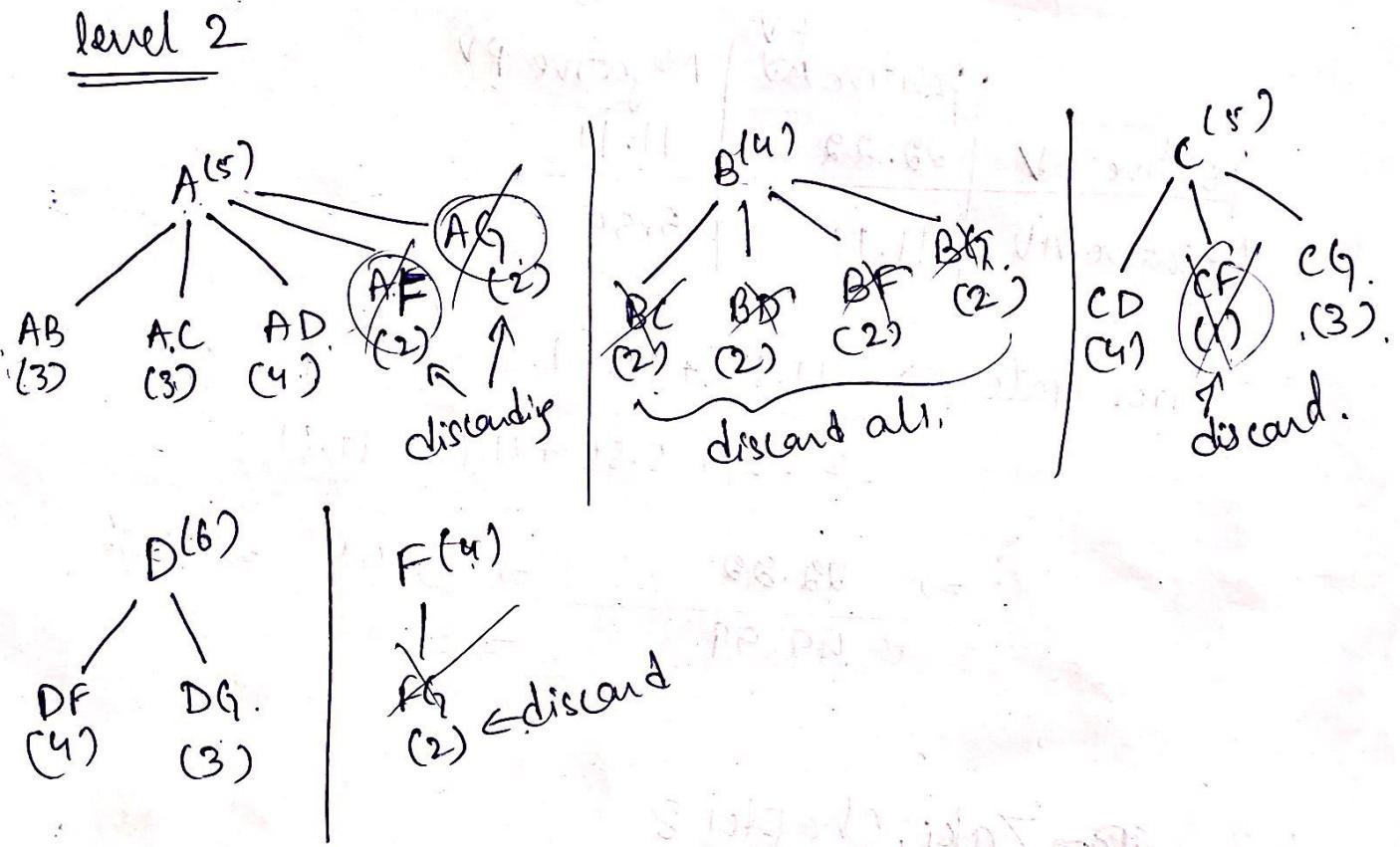
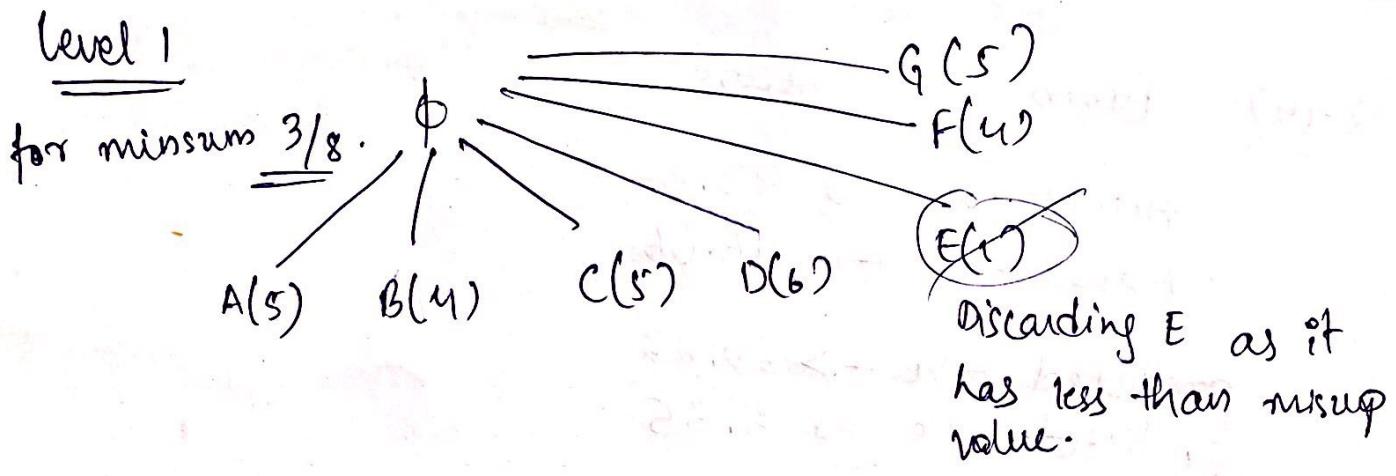
$$\Rightarrow \frac{22.22}{49.99} \Rightarrow 0.44 = 44\%$$

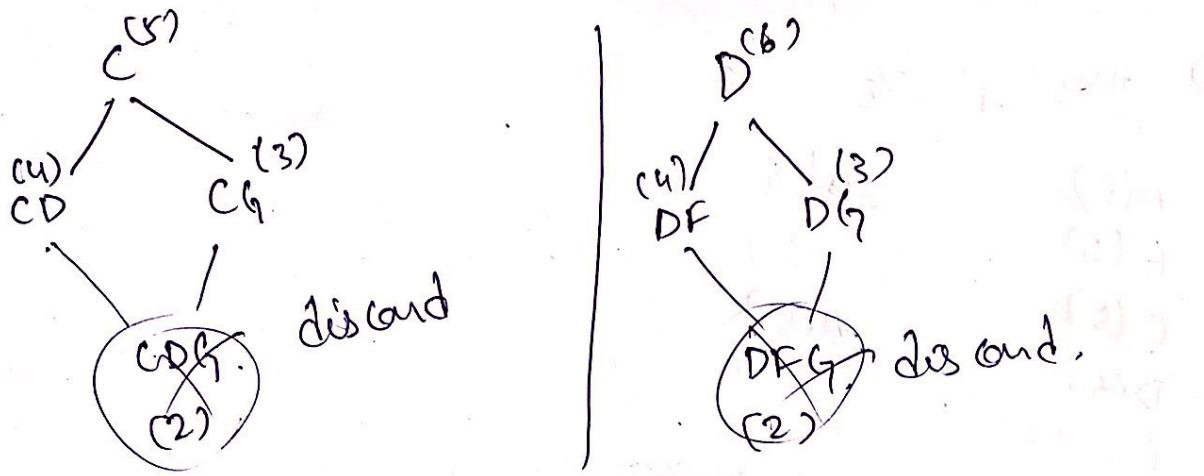
1.3 ~~Zaki~~, Chapter 8

Ex 1)

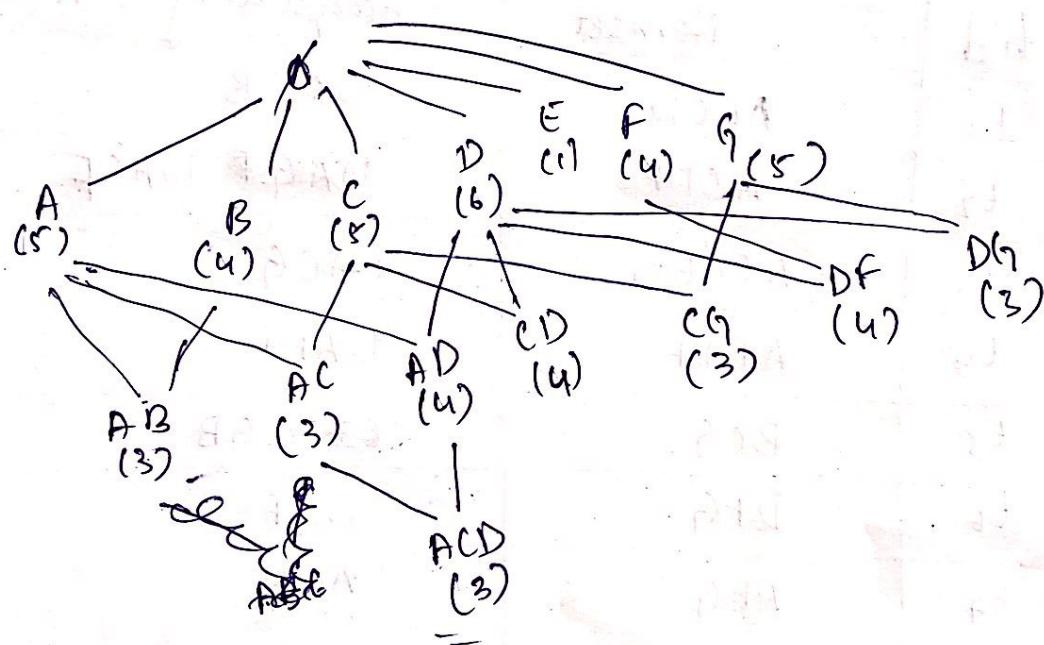
a)	tid	itemset
	t ₁	ABCD
	t ₂	ACDF
	t ₃	ACDEG
	t ₄	ABDF
	t ₅	BCG
	t ₆	DFG
	t ₇	ABG
	t ₈	CDFG.

	A	B	C	D	E	F	G
t ₁	1	1	1	1			
t ₂	1		1	1	1		
t ₃	1		1	1	1	1	
t ₄	1	1		1	1	1	
t ₅		1	1				1
t ₆				1	1	1	
t ₇	1	1		1	1	1	
t ₈	1	1	1	1	1	1	1





Overall would be

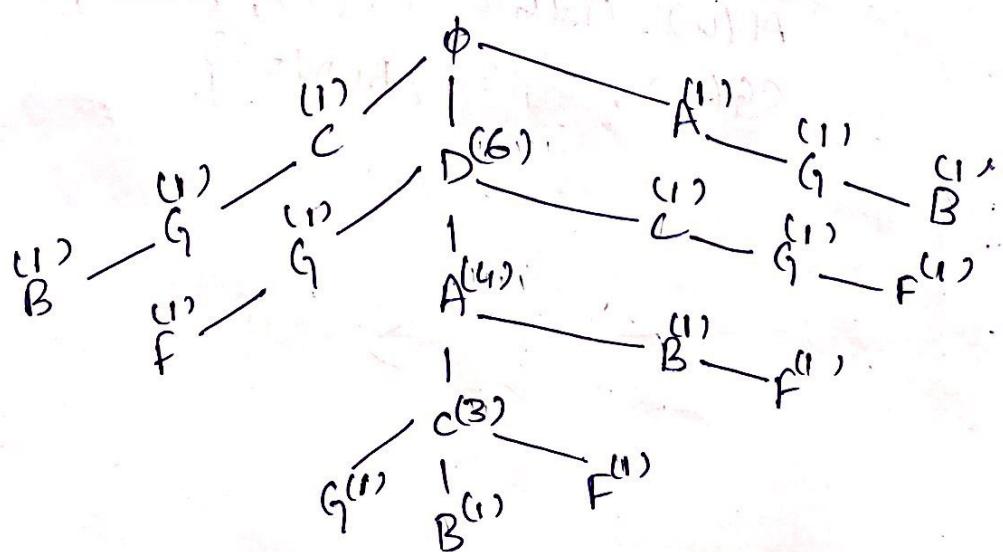


Fragment items = $\{D(6), A(5), C(5), G(5), B(4), F(4),$
 $AD(4), CD(4), DF(4), AB(3), AC(3),$
 $CG(3), DG(3), ACD(3)\}.$

1-(b) $\text{minsup} = \frac{1}{8}$.

	<u>sorted</u>
A(5)	
B(4)	D(6)
C(5)	A(5)
D(6)	C(5)
E(1)X	G(5)
F(4)	B(4)
G(5)	F(4)

tid	itemset	<u>frequently sorted</u>
t ₁	ABCD	DACB
t ₂	ACDF	DA G CF
t ₃	ACDEG	DACG
t ₄	ABDF	DABF
t ₅	BCG	CBG
t ₆	DFG	DGF
t ₇	ABG	AGB
t ₈	CDFG	DCGF



projection :-

$$(1) R_A \Rightarrow D^{(4)} \Rightarrow \boxed{\{AD^{(4)}\}}$$

$$\begin{array}{c} \phi \\ | \\ D^{(4)} \end{array}$$

$$(2) R_B \Rightarrow DAC, DA, CG, AG.$$

$$DAC = \{D(2), A(3), C(2), DA(2)\}$$

~~$$DA = \{D(2), A(3), DA(2)\}$$~~

$$DA = \{D(2), A(3), DA(2)\}$$

$$CG = \{C(2), G(2), CG(1)\}$$

$$AG = \{A(3), G(2), AG(1)\}$$

$$R_B = \{D(2), A(3), C(2), DA(2), G(2)\}$$

$$= \{BD(2), BA(3), BC(2), BDA(2), BG(2)\}$$

$$(3) R_C \Rightarrow DA, D$$

$$DA = \{D(2), A(1), DA(1)\}$$

$$D = \{CD(2)\}$$

$$(1) R_D \Rightarrow \text{None}$$

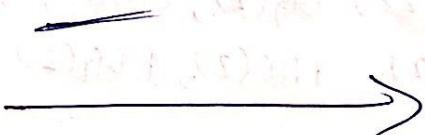
$$(3) R_C \Rightarrow DA, D.$$

$$DA = \{D(4), A(3), DA(3)\}$$

$$(5) R_F \Rightarrow$$

$$= \{CD(4), CA(3), CDA(3)\}$$

$$(4) R_D \Rightarrow \text{None}$$



(5) $R_F = DAC, DAB, DG, DCg$.

$DAC = \{D(4), A(2), C(2), DA(2), DC(2), AC(1)\}$

$DAB = \{D(4), A(2), B(1), DA(2), DB(1), AB(1)\}$

$DG = \{D(4), G(2), DG(2)\}$

$DCg = \{D(4), C(2), G(2), DC(2), CG(1), DG(2)\}$

$R_F = \{D(4), A(2), C(2), DA(2), DC(2), DG(2)\}$

$\Rightarrow \{FD(4), \cancel{FA(2)}, FC(2), FDA(2), FDC(2), FDG(2)\}$

(6) $R_G = DAC, C, D, A, DC$

$DAC = \{D(3), A(2), C(3), DA(1), DC(2), AC(1)\}$

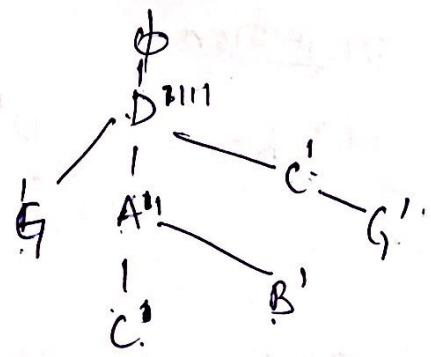
$\Leftrightarrow DC = \{D(3), C(3), DC(2)\}$

$R_G = \{D(3), A(2), C(3), DC(2)\}$

$\Rightarrow \{GD(3), GA(2), GC(3), GDC(2)\}$

The frequent itemset $\Rightarrow 25$ items

$\{D(6), A(5), C(5), G(5), B(4), F(4), AD(4), BD(2), BA(3), BC(2), BDA(2), BG(2), CD(4), CA(3), CDA(3), FD(4), FA(2), FC(2), FDA(2), FDC(2), FDG(2), GD(3), GA(2), GC(3), GDC(2)\}$



Q4)

tid	itemset
t ₁	ACD
t ₂	BCE
t ₃	ABCDE
t ₄	BDE
t ₅	ABCE
t ₆	ABCD

Rules that one can generate from the set ABE

$$\Rightarrow \text{let } a = \{ABE\}$$

calculating all the subsets of a.

$\therefore a = \{AB(3), AE(2), BE(1), A(1), B(1), E(1), ABE(2), \emptyset(6)\}$.
for each subset, generating a rule.

$$* \{AB\} \rightarrow \{E\}$$

$$\text{confidence (C)} = \frac{\text{supp}(ABE)}{\text{supp}(AB)} \Rightarrow \frac{2}{3} = 0.66$$

$$* \{AE\} \rightarrow \{B\}$$

$$C = \frac{\text{supp}(ABE)}{\text{supp}(AE)} \Rightarrow \frac{2}{2} = 1.$$

$$* \{BE\} \rightarrow \{A\}$$

$$C = \frac{\text{supp}(ABE)}{\text{supp}(BE)} \Rightarrow \frac{2}{4} = \frac{1}{2} = 0.5.$$

$$* \{ABE\} \rightarrow \{\emptyset\}$$

$$C = \frac{2}{2} = 1$$



* $\{A\} \rightarrow \{BCE\}$

$$C = \frac{\text{supp}(ABCE)}{\text{supp}(A)} = \frac{2}{4} = 0.5$$

* $\{B\} \rightarrow \{ACE\}$

$$C = \frac{\text{supp}(ABCE)}{\text{supp}(B)} = \frac{2}{5} = 0.4. \quad \left| \begin{array}{l} \{B\} \rightarrow \{ABCE\} \\ C = \frac{2}{6} = \frac{1}{3} \end{array} \right.$$

* $\{E\} \rightarrow \{ABC\}$

$$C = \frac{\text{supp}(ABCE)}{\text{supp}(E)} = \frac{2}{4} = 0.5.$$

Q 6)

6-a) Number of itemsets $\Rightarrow 2^k - 1 = 2^{11} - 1 = 2047$.

6-b) More than or equal to support of x.
(option iv).

1.4 Multiclass classification

The multiclass.Rmd's confusion matrix created by 'R'.
is as follows :-

predicted value

		setosa	versicolor	virginica
Actual value {	setosa	8	0	0
	versicolor	0	10	2
	virginica	0	1	9

→ for setosa, the binary confusion matrix is,

	setosa	vericolor + virginica
setosa	8	0.
vericolor + virginica	0.	22

$$\text{sensitivity} \Rightarrow \frac{TP}{TP+FN} = \frac{8}{8+0} = 1.0$$

$$\text{specificity} \Rightarrow \frac{TN}{TN+FP} = \frac{22}{22+0} = 1.0$$

$$\text{precision} = \frac{TP}{TP+FP} = \frac{8}{8+0} = 1.0$$

→ for vericolor

	vericolor	PV
	vericolor	setosa + virginica
setosa	10.	2.
setosa + virginica	1	17.

$$\text{sensitivity} \Rightarrow \frac{10}{10+2} = \frac{10}{12} = 0.83$$

$$\text{specificity} \Rightarrow \frac{17}{18} = 0.94$$

$$\text{precision} \Rightarrow \frac{10}{11} = 0.91$$

→ for virginica,

	P	V
virginica.	9	1
setosa + vericolor	2	18

$$\text{sensitivity} \Rightarrow \frac{9}{10} = 0.90$$

$$\text{specificity} \Rightarrow \frac{18}{20} = 0.90.$$

$$\text{precision} \Rightarrow \frac{9}{11} = 0.82$$