

11/2/18

Data Mining

Homework-2

A20406657

Avinash Vellamuri

Problem 1.2:

2) a) Compute the gini index for the overall collection of training examples.

20 observation $\begin{cases} \rightarrow 10 \text{ are class 0} \\ \rightarrow 10 \text{ are class 1} \end{cases}$

$$\begin{aligned} \text{gini index} &= 1 - \left[\left(\frac{10}{20} \right)^2 + \left(\frac{10}{20} \right)^2 \right] = 1 - [0.25 + 0.25] \\ &= 1 - 0.5 = 0.5 \end{aligned}$$

$$\boxed{\text{gini index} = 0.5}$$

b) gini index for the customer ID attribute.

Since the customer ID's are unique (ie) no relation between the customer ID values.

\therefore The gini index for Customer ID is Zero

c) Gini index for Gender attribute.

we have 10 observations as Male and 10 observations as Female

	M	F
C ₀	6	4
C ₁	4	6

$$\text{Gini-Index (Male)} = 1 - \left[\left(\frac{6}{10} \right)^2 + \left(\frac{4}{10} \right)^2 \right] = 0.48$$

$$\text{Gini-Index (Female)} = 1 - \left[\left(\frac{4}{10} \right)^2 + \left(\frac{6}{10} \right)^2 \right] = 0.48$$

$$\text{Overall gini Index} = \frac{10}{20} \times 0.48 + \frac{10}{20} \times 0.48 = 0.48$$

d) Gini index for car type using multiway split.

	Car type		
	Family	Sports	Luxury
Co	1	8	1
C ₁	3	0	7

$$\text{Gini}(\text{Family}) = 1 - \left[\left(\frac{1}{4} \right)^2 + \left(\frac{3}{4} \right)^2 \right] = 0.375$$

$$\text{Gini}(\text{Sports}) = 1 - \left[\left(\frac{8}{8} \right)^2 + \left(\frac{0}{8} \right)^2 \right] = 0$$

$$\text{Gini}(\text{Luxury}) = 1 - \left[\left(\frac{1}{8} \right)^2 + \left(\frac{7}{8} \right)^2 \right] = 0.21875$$

$$\begin{aligned} \text{Overall Gini} &= \frac{4}{20} \times 0.375 + \frac{8}{20} \times 0 + \frac{8}{20} \times 0.21875 \\ &= 0.075 + 0.0875 = \underline{0.1625} \end{aligned}$$

e) Gini index for Shirt size using multiway split.

	Shirt size			
	Small	Medium	Large	Extra large
Co	3	3	2	2
C ₁	2	4	2	2

$$\text{Gini}(\text{Small}) = 1 - \left[\left(\frac{3}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right] = 0.48$$

$$\text{Gini}(\text{Medium}) = 1 - \left[\left(\frac{3}{7} \right)^2 + \left(\frac{4}{7} \right)^2 \right] = 0.4898$$

$$\text{Gini}(\text{Large}) = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = 0.5$$

$$\text{Gini}(\text{Extra large}) = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = 0.5$$

Overall Gini Index of shirt size

$$= \frac{5}{20} \times 0.48 + \frac{7}{20} \times 0.4898 + \frac{4}{20} \times 0.5 + \frac{4}{20} \times 0.5$$

$$= 0.49143$$

p) which among Gender, Car Type, Shirt size is better?
Car type is better since it has the lowest gini index.

q) why Customer ID should not be used as the attribute test condition?

Customer ID is unique. Each customer is assigned with a unique ID called the Customer ID.
We can't predict the customer ID attribute because new customers are given unique ID.

3)

a) Entropy of the training examples with respect to the positive class?

$$\text{Entropy (positive class)} = - \left[\frac{4}{9} \log_2 \frac{4}{9} + \frac{5}{9} \log_2 \frac{5}{9} \right]$$

$$= 0.9911$$

	<u>a₁</u>	
	+	-
T	3	1
F	1	4

	<u>a₂</u>	
	+	-
T	2	3
F	2	2

$$\text{Entropy}(a_1) = \frac{4}{9} \left[-\left(\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4} \right) \right] + \frac{5}{9} \left[-\left(\frac{1}{5} \log_2 \frac{1}{5} + \frac{4}{5} \log_2 \frac{4}{5} \right) \right]$$

$$= 0.7616$$

$$\text{Entropy}(a_2) = \frac{4}{9} \left[-\left(\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4} \right) \right] + \frac{5}{9} \left[-\left(\frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5} \right) \right]$$

$$= 0.9839$$

$$\text{Information Gain}(a_1) = 0.9911 - 0.7616 = 0.2294$$

$$\text{Information gain}(a_2) = 0.9911 - 0.9839 = 0.0072$$

c) For a_3 find information gain for every possible split.

	1	3	4	5	6	7
split →	2	3.5	4.5	5.5	6.5	7.5
	<= >	<= >	<= >	<= >	<= >	<= >
+	1 3	1 3	2 2	2 2	3 1	4 0
-	0 5	1 4	1 4	3 2	3 2	4 1

Split 2

$$\frac{1}{9} \left[-\left(\frac{1}{1} \log_2 1 + \frac{0}{1} \log_2 0 \right) \right] + \frac{8}{9} \left[-\left(\frac{3}{8} \log_2 \frac{3}{8} + \frac{5}{8} \log_2 \frac{5}{8} \right) \right]$$

$$\text{Entropy} = 0.8486$$

$$\text{Information gain} = 0.9911 - 0.8486 = 0.1425$$

Split 35

$$\frac{2}{9} \left[-\left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right) \right] + \frac{7}{9} \left[-\left(\frac{3}{7} \log_2 \frac{3}{7} + \frac{4}{7} \log_2 \frac{4}{7} \right) \right]$$

$$\text{Entropy} = 0.9885$$

$$\text{Information gain} = 0.9911 - 0.9885 = 0.0026$$

Split 4.5:-

$$\frac{3}{9} \left[- \left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right) \right] + \frac{6}{9} \left[- \left(\frac{2}{6} \log_2 \frac{2}{6} + \frac{4}{6} \log_2 \frac{4}{6} \right) \right]$$

$$\text{Entropy} = 0.9183$$

$$\text{Info gain} = 0.9911 - 0.9183 = 0.0728$$

Split 5.5:-

$$\frac{5}{9} \left[- \left(\frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5} \right) \right] + \frac{4}{9} \left[- \left(\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4} \right) \right]$$

$$\text{Entropy} = 0.9839$$

$$\text{Info gain} = 0.9911 - 0.9839 = 0.0072$$

Split 6.5:-

$$\frac{6}{9} \left[- \left(\frac{3}{6} \log_2 \frac{3}{6} + \frac{3}{6} \log_2 \frac{3}{6} \right) \right] + \frac{3}{9} \left[- \left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3} \right) \right]$$

$$\text{Entropy} = 0.9728$$

$$\text{Info gain} = 0.9911 - 0.9728 = 0.0183$$

Split 7.5:-

$$\frac{8}{9} \left[- \left(\frac{4}{8} \log_2 \frac{4}{8} + \frac{4}{8} \log_2 \frac{4}{8} \right) \right] + \frac{1}{9} \left[- (0 \log_2 0 + 1 \log_2 1) \right]$$

$$\text{Entropy} = 0.8889$$

$$\text{Info gain} = 0.9911 - 0.8889 = 0.1022$$

Best split for a_3 is at position split 2

Here the information gain is maximum.

d) From the information gains of attributes a_1, a_2, a_3 .

$$a_1 = 0.2294; \quad a_2 = 0.0072; \quad a_{3_{\max}} = 0.1427$$

a_1 is the best attribute to split as it has high information gain.

e) q1

$$\text{Error rate}(T) = 1 - \left[\frac{3}{4}, \frac{1}{4} \right]_{\max} = 1 - \frac{3}{4} = \frac{1}{4}$$

$$\text{Error rate}(F) = 1 - \text{Max} \left[\frac{1}{5}, \frac{4}{5} \right] = 1 - \frac{4}{5} = \frac{1}{5}$$

$$\text{Error}(a_1) = \frac{4}{9} \left(\frac{1}{4} \right) + \frac{5}{9} \left(\frac{1}{5} \right) = \underline{\underline{\frac{2}{9}}}$$

a2

$$\text{Error rate}(T) = 1 - \left[\text{Max} \left[\frac{2}{4}, \frac{2}{4} \right] \right] = 1 - \frac{1}{2} = \frac{1}{2}$$

$$\text{Error rate}(F) = 1 - \text{Max} \left[\frac{3}{5}, \frac{2}{5} \right] = 1 - \frac{3}{5} = \frac{2}{5}$$

$$\text{Error}(a_2) = \frac{4}{9} \left(\frac{1}{2} \right) + \frac{5}{9} \left(\frac{2}{5} \right) = \frac{4}{9}$$

a_1 is best than a_2 as a_2 has higher error rate.

d) what is the best split according to the Gini index?

q1

$$\text{Gini}(+ve) = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 0.375$$

$$\text{Gini}(-ve) = 1 - \left[\left(\frac{1}{5} \right)^2 + \left(\frac{4}{5} \right)^2 \right] = 0.32$$

$$\text{Gini}(a_1) = \frac{4}{9} \times 0.375 + \frac{5}{9} \times 0.32 = 0.344$$

a2

$$\text{Gini}(a_2) = \frac{4}{9} \left[1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] \right] + \frac{5}{9} \left[1 - \left[\left(\frac{3}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right] \right]$$

$$= 0.4889$$

$\therefore \text{Gini}(a_1) < \text{Gini}(a_2)$. a_1 is a better split.

5) a) Information gain when splitting on A and B.

+ve 4 out of 10, -ve 6 out of 10.

(A)	+	-
T	4	3
F	0	3

$$E_A(+ve) = - \left[\frac{4}{4} \log_2 \frac{4}{4} + \frac{0}{4} \log_2 \frac{0}{4} \right]$$

(A)	T	F
+	4	0
-	3	3

(B)	T	F
+	3	1
-	1	5

$$E_{\text{tot}} = - \left[\frac{4}{10} \log_2 \frac{4}{10} + \frac{6}{10} \log_2 \frac{6}{10} \right]$$

$$E_{\text{tot}} = 0.9710$$

Split on A:-

$$E_T = - \left[\frac{4}{7} \log_2 \frac{4}{7} + \frac{3}{7} \log_2 \frac{3}{7} \right] = 0.9852$$

$$E_F = - \left[\frac{0}{3} \log_2 \frac{0}{3} + \frac{3}{3} \log_2 \frac{3}{3} \right] = 0$$

$$\Delta = E_{\text{tot}} - \left[\frac{7}{10} \times 0.9852 \right] = 0.9710 - 0.68964$$

$$\Delta_{\text{info}} = 0.2813$$

split on B:-

$$E_T = - \left[\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4} \right] = 0.8113$$

$$E_F = - \left[\frac{1}{6} \log_2 \frac{1}{6} + \frac{5}{6} \log_2 \frac{5}{6} \right] = 0.65$$

$$\Delta = 0.9710 - \left[\frac{4}{10} \times 0.8113 + \frac{6}{10} \times 0.65 \right]$$

$$\Delta = 0.2569$$

Information gain for 'A' is higher. So it is better to split at A.

$$b) G_{\text{tot}} = 1 - \left[\left(\frac{4}{10} \right)^2 + \left(\frac{6}{10} \right)^2 \right] = 0.48$$

Splitting at A:-

$$G_T = 1 - \left[\left(\frac{4}{7} \right)^2 + \left(\frac{3}{7} \right)^2 \right] = 0.4898$$

$$G_F = 1 - \left[\left(\frac{0}{3} \right)^2 + \left(\frac{3}{3} \right)^2 \right] = 0$$

$$\Delta = 0.48 - \left[\frac{7}{10} \times 0.4898 \right] = 0.1371$$

Splitting at B:-

$$G_T = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 0.3750$$

$$G_F = 1 - \left[\left(\frac{1}{6} \right)^2 + \left(\frac{5}{6} \right)^2 \right] = 0.2778$$

$$\Delta = 0.48 - \left[\frac{4}{10} \times 0.3750 + \frac{6}{10} \times 0.2778 \right]$$
$$= 0.1633$$

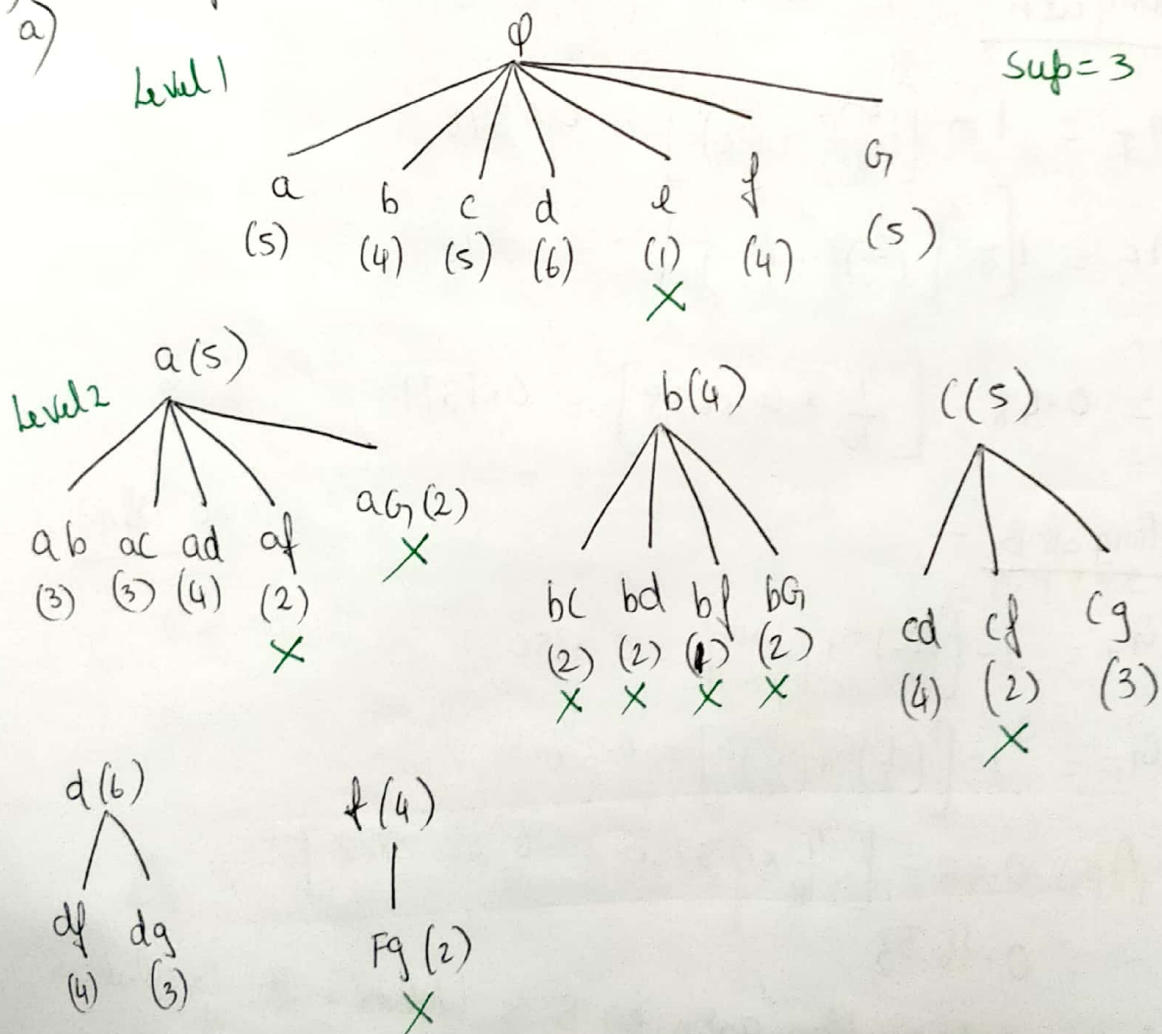
Since information gain for B is higher. It is better to split at attribute B. Algorithm picks "B"

c) Yes, An attribute selection can vary based upon gini index or entropy " Δ " calculation we can see from the results of a) and b) of question 5 that both the method picked different attributes.

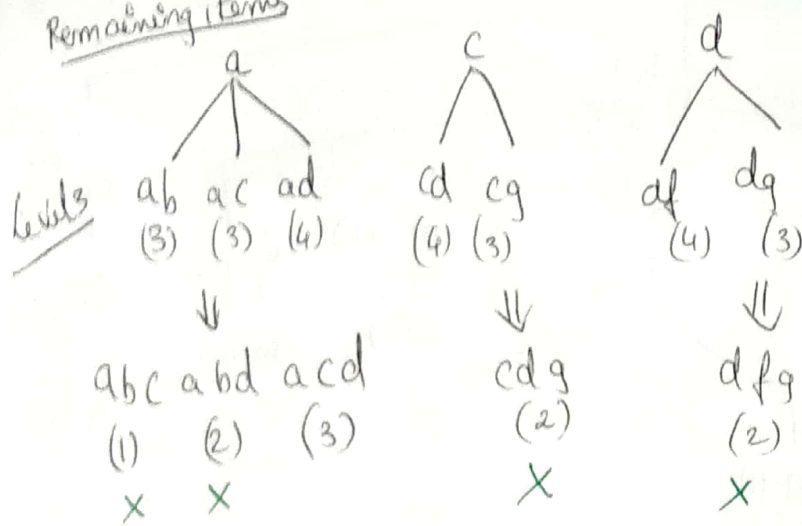
Because the gain calculated " Δ " are scaled differently. Entropy follows a long scale while gini index follows a polynomial scale.

Problem 1.4 (Zaki)

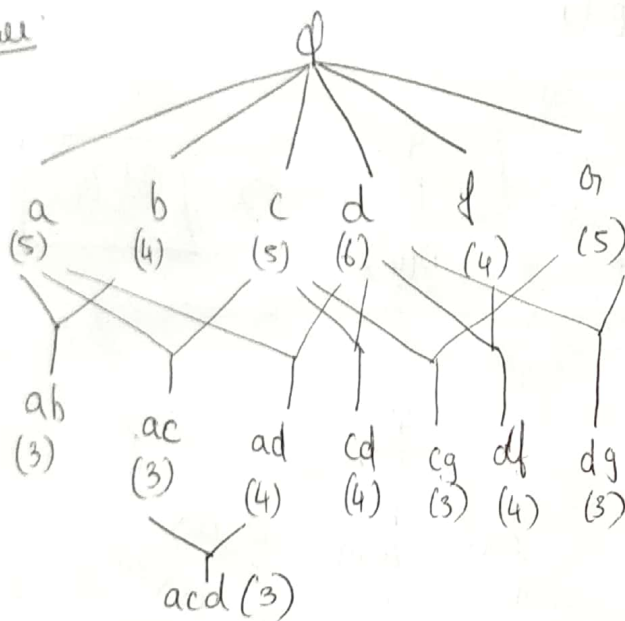
1) $\text{Minsup} = 3/8$, frequent patterns using apriori



Remaining items



Final tree

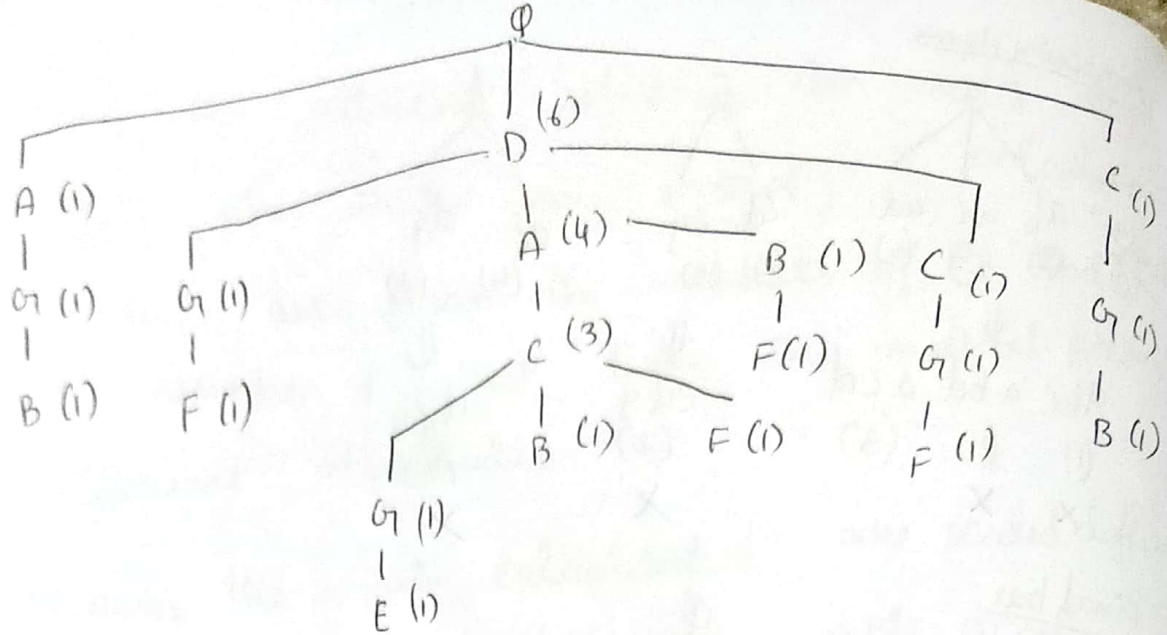


Frequent Item Sets = $\{d(6), g(5), a(5), c(5), b(4), f(4), df(4), ad(4), cd(4), ab(3), cg(3), dg(3), ac(3), acd(3)\}$
 "14 items"

1b) $D(6) \quad A(5) \quad C(5) \quad G(5) \quad B(4) \quad F(4) \quad E(1)$
 X

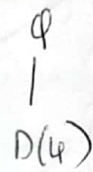
$ABCD \rightarrow DACB$
 $ACDF \rightarrow DACF$
 $A(DEG) \rightarrow DACGE$
 $ABDF \rightarrow DABF$
 $BCG \rightarrow CGB$
 $DFG \rightarrow DGF$
 $ABG \rightarrow AGB$
 $CD FG \rightarrow DCFG$

$D(6), A(5), C(5), G(5)$
 $B(4), F(4)$



Projection trees:-

R_A $D, \text{count}(A)=4$



\Rightarrow Ad(4)

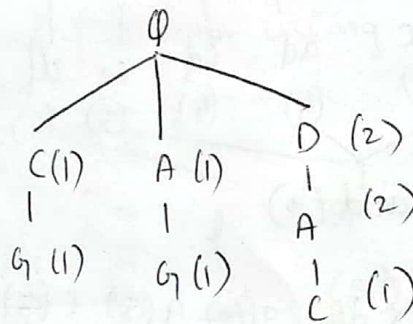
R_B

$AG, \text{count}(B)=1$

$DAC, \text{count}(B)=1$

$DA, \text{count}(B)=1$

$CG, \text{count}(B)=1$



$\{D, A, C\} \Rightarrow \{D(2), A(2), C(1), DA(2), DC(1), AC(1)\}$

$\{A, G\} \Rightarrow \{A(1), G(1), AG(1)\}$

$\{C, G\} \Rightarrow \{C(1), G(1), CG(1)\}$

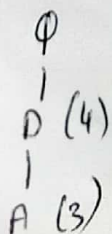
$\Rightarrow \{D(2), A(3), C(2), G(2), DA(2)\}$

$BD(2), BA(3), BC(2), BG(2), BDA(2)$

R_C

$DA, \text{count}(C)=3$

$D, \text{count}(C)=1$

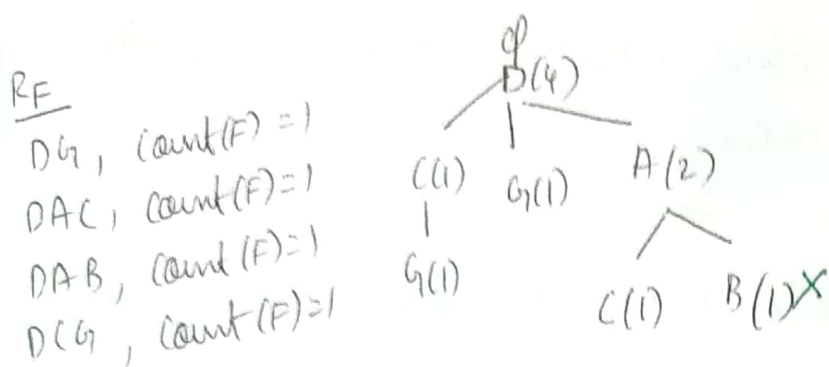


$\{D(4), A(3), DA(3)\}$

$CD(4), CA(3), CDA(3)$

$R_D \Rightarrow \text{None}$

R_E
 $DACG, \text{count}(E) = 1 \Rightarrow \text{None}$ all have count 1.



$\{D, A\} \Rightarrow \{D(4), A(2), DA(2)\}$

$\boxed{FD(4), FA(2), FDA(2)}$

$\{D, G\} \Rightarrow \{D(4), G(1), DG(1)\}$

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

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

$\boxed{FG(2), FDG(2), FDC(2), FC(2)}$

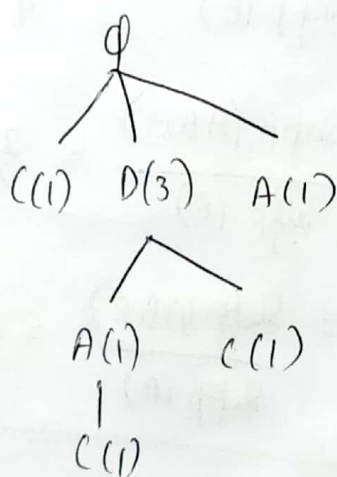
R_G

$A, \text{count}(G) = 1$

$D, \text{count}(G) = 1$

$DAC, \text{count}(G) = 1$

$DC, \text{count}(G) = 1$



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

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

Total no of frequent item sets $\Rightarrow 26$ sets

D(6), A(5), C(5), G(5), B(4), F(4), Ad(4), BD(2), BA(3),
B(2), BG(2), BDA(2), CD(4), CA(3), CDA(3), FD(4),
FA(2), FDA(2), FG(2), FDG(2), FDC(2), FC(2), GA(2),
GD(3), GDC(2).

4) $ABE \Rightarrow \{A(4), B(5), E(4), AB(3), AE(2), BE(4)\}$

$$C(BE \rightarrow A) = \frac{\text{Supp}(ABE)}{\text{Supp}(BE)} = \frac{2}{4} = 0.5$$

$$C(AE \rightarrow B) = \frac{\text{Supp}(ABE)}{\text{Supp}(AE)} = \frac{2}{2} = 1$$

$$C(AB \rightarrow E) = \frac{\text{Supp}(ABE)}{\text{Supp}(AB)} = \frac{2}{3} = 0.66$$

$$C(E \rightarrow AB) = \frac{\text{Supp}(ABE)}{\text{Supp}(E)} = \frac{2}{4} = 0.5$$

$$C(B \rightarrow AE) = \frac{\text{Supp}(ABE)}{\text{Supp}(B)} = \frac{2}{5} = 0.4$$

$$C(A \rightarrow BE) = \frac{\text{Supp}(ABE)}{\text{Supp}(A)} = \frac{2}{4} = 0.5$$

6) here $k=11$

a) No of itemsets $= 2^k - 1 = 2^{11} - 1 = 2047$

b) (iv) More than or equal to support of X