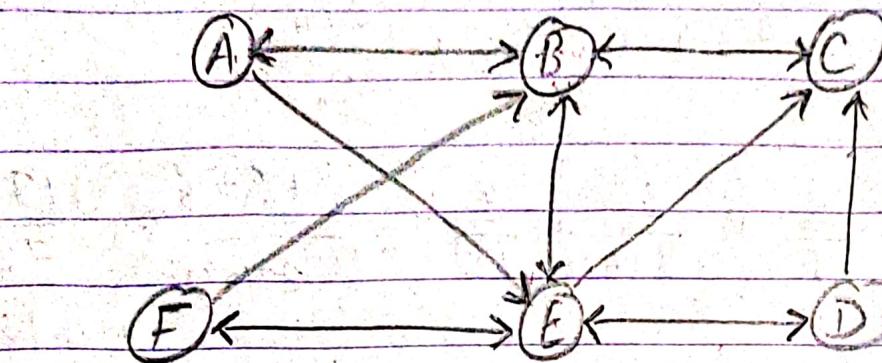


* Page Rank Algorithm



upto round 2

$$\text{Formula: } PR(U) = \sum_{V \in B_U} \frac{PR(V)}{L(V)}$$

$$\begin{aligned} \text{Initially, } PR(A) &= PR(B) = PR(C) = PR(D) = PR(E) \\ &= PR(F) = \frac{1}{n} = \frac{1}{6} \end{aligned}$$

where, n = number of nodes.

Iteration : 1

$$PR(A) = \frac{PR(B)}{L(B)} = \frac{\frac{1}{6}}{3} = \frac{1}{18} \approx 0.055$$

$$PR(B) = \frac{PR(A)}{L(A)} + \frac{PR(C)}{L(C)} + \frac{PR(E)}{L(E)} + \frac{PR(F)}{L(F)}$$

$$= \frac{\frac{1}{6}}{2} + \frac{\frac{1}{6}}{1} + \frac{\frac{1}{6}}{4} + \frac{\frac{1}{6}}{2}$$

$$= \frac{1}{12} + \frac{1}{6} + \frac{1}{24} + \frac{1}{12} = \frac{2+4+1+2}{24} = \frac{9}{24}$$

$$= 0.375$$

$$PR(C) = \frac{PR(B)}{L(B)} + \frac{PR(E)}{L(E)} + \frac{PR(D)}{L(D)}$$

$$= \frac{1/6}{3} + \frac{1/6}{4} + \frac{1/6}{2}$$

$$= \frac{1}{18} + \frac{1}{24} + \frac{1}{12} = \frac{4+3+6}{72} = \frac{13}{72} = 0.18$$

$$PR(D) = \frac{PR(E)}{L(E)}$$

$$= \frac{1/6}{4} = \frac{1}{24} = 0.041$$

$$PR(E) = \frac{PR(A)}{L(A)} + \frac{PR(B)}{L(B)} + \frac{PR(D)}{L(D)} + \frac{PR(F)}{L(F)}$$

$$= \frac{1/6}{2} + \frac{1/6}{3} + \frac{1/6}{2} + \frac{1/6}{2}$$

$$= \frac{1}{12} + \frac{1}{18} + \frac{1}{12} + \frac{1}{12}$$

$$= \frac{3+2+3+2}{36} = \frac{11}{36} = 0.305$$

$$PR(F) = \frac{PR(E)}{L(E)}$$

$$= \frac{1/6}{4}$$

$$= \frac{1}{24} = 0.0416$$

Iteration: 2

$$PR(A) = \frac{PR(B)}{L(B)} = \frac{0.375}{3} = 0.125$$

$$\begin{aligned} PR(B) &= \frac{PR(A)}{L(A)} + \frac{PR(C)}{L(C)} + \frac{PR(E)}{L(E)} + \frac{PR(F)}{L(F)} \\ &= \frac{0.055}{2} + \frac{0.18}{1} + \frac{0.305}{4} + \frac{0.0416}{2} \\ &= ? \end{aligned}$$

$$\begin{aligned} PR(C) &= \frac{PR(B)}{L(B)} + \frac{PR(E)}{L(E)} + \frac{PR(D)}{L(D)} \\ &= \frac{0.375}{3} + \frac{0.305}{4} + \frac{0.041}{2} \\ &= ? \end{aligned}$$

$$PR(D) = \frac{PR(E)}{L(E)} = \frac{0.305}{4} = 0.076$$

$$\begin{aligned} PR(E) &= \frac{PR(A)}{L(A)} + \frac{PR(B)}{L(B)} + \frac{PR(D)}{L(D)} + \frac{PR(F)}{L(F)} \\ &= \frac{0.055}{2} + \frac{0.375}{3} + \frac{0.041}{2} + \frac{0.0416}{2} = ? \end{aligned}$$

$$PR(F) = \frac{PR(E)}{L(E)} = \frac{0.305}{4} = 0.076$$

* HITS Algorithm:

'Hypertext Induced Topic Search'

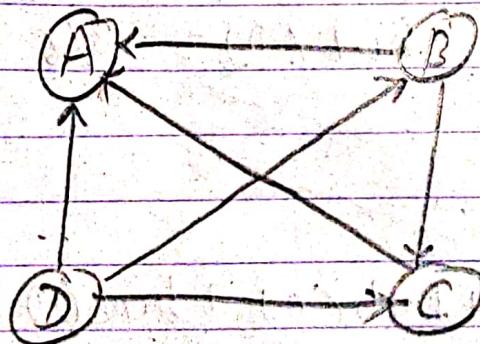
Formula:

Authority Update Rule:

- Each node's authority score is the sum of hub scores of each node that points to it. (Inbound links)

Hub Update Rule:

- Each node's hub score is the sum of authority scores of each node that it points to. (outbound links)



Initially;

$$\begin{aligned} \text{auth}(A) &= \text{auth}(B) = \text{auth}(C) = \text{auth}(D) = 1 \\ \text{hub}(A) &= \text{hub}(B) = \text{hub}(C) = \text{hub}(D) = 1 \end{aligned}$$

Iteration: 1

$$\text{auth}(A) = \text{hub}(B) + \text{hub}(C) + \text{hub}(D) = 3$$

$$\text{auth}(B) = \text{hub}(D) = 1$$

$$\text{auth}(C) = \text{hub}(B) + \text{hub}(D) = 2$$

$$\text{auth}(D) = 0$$

$$\text{hub}(A) = \cancel{\text{auth}}(A) = 0$$

$$\text{hub}(B) = \text{auth}(A) + \text{auth}(C) = 2$$

$$\text{hub}(C) = \text{auth}(A) = 1$$

$$\text{hub}(D) = \text{auth}(A) + \text{auth}(B) + \text{auth}(C) = 3$$

After applying normalization:

$$\text{auth}(A) = 3/6$$

$$\text{hub}(A) = 0$$

$$\text{auth}(B) = 1/6$$

$$\text{hub}(B) = 2/6$$

$$\text{auth}(C) = 2/6$$

$$\text{hub}(C) = 1/6$$

$$\text{auth}(D) = 0$$

$$\text{hub}(D) = 3/6$$

Iteration: 2

$$\text{auth}(A) = \text{hub}(B) + \text{hub}(C) + \text{hub}(D)$$

$$= \frac{2}{6} + \frac{1}{6} + \frac{3}{6} = \frac{6}{6} = 1$$

$$\text{auth}(B) = \text{hub}(D) = 3/6$$

$$\text{auth}(C) = \text{hub}(B) + \text{hub}(D) = \frac{2}{6} + \frac{3}{6} = \frac{5}{6}$$

$$\text{auth}(D) = 0$$

$$\text{hub}(A) = 0$$

$$\text{hub}(B) = \text{auth}(A) + \text{auth}(C)$$

$$= \frac{3}{6} + \frac{2}{6} = \frac{5}{6}$$

$$\text{hub}(C) = \text{auth}(A) = \frac{3}{6}$$

$$\text{hub}(D) = \text{auth}(A) + \text{auth}(B) + \text{auth}(C)$$

$$= \frac{3}{6} + \frac{1}{6} + \frac{2}{6} = \frac{6}{6} = 1$$

After applying normalization:

$$\text{auth}(A) = \frac{1}{1+3/6+5/6+0} = \frac{6}{14} \quad \text{(u)}$$

$$\text{auth}(B) = \frac{3}{6} \times \frac{6}{14} = \frac{3}{14}$$

$$\text{auth}(C) = \frac{5}{6} \times \frac{6}{14} = \frac{5}{14}$$

$$\text{auth}(D) = 0$$

No^o, $\text{hub}(A) = 0$

$$\text{hub}(B) = \frac{5}{6} \times \frac{6}{14} = \frac{5}{14}$$

$$\text{hub}(C) = \frac{3}{6} \times \frac{6}{14} = \frac{3}{14}$$

$$\text{hub}(D) = \frac{6}{14}$$

* Collaborative Filtering Recommender Systems

Similarity functions!

$$1. \text{ Pearson } C(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}}$$

$$2. \text{ Cosine } C(x, y) = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}}$$

(User-based)

Example 1; Suppose the ratings of five users U_1, U_2, U_3, U_4 and U_5 are drawn for six items denoted by I_1, I_2, I_3, I_4, I_5 and I_6 . Each rating is drawn from the range basis of $\{1, \dots, 7\}$.

Find the predictions of user U_3 for items I_1 and I_6 on the basis of the ratings in the table below. Use Pearson Correlation Coefficient to find similarity between users and consider top-2 closest users.

Row-wise

Row-wise

User/Item	I1	I2	I3	I4	I5	I6
U1	7	6	7	4	5	4
U2	6	7	?	4	3	4
U3	?	3	3	1	1	?
U4	1	2	2	3	3	4
U5	1	?	1	2	3	3

Soln'

The first step is to compute the similarity between user U3 and all the other users.

$$\text{Pearson } (U1, U3) = \frac{(6-5.5)(3-2) + (7-5.5)(3-2)}{+(4-5.5)(1-2) + (5-5.5)(1)}$$

$$\sqrt{1.5^2 + 1.5^2 + (-1.5)^2 + (-0.5)^2}$$

$$\sqrt{1^2 + 1^2 + (-1)^2 + (-1)^2}$$

$$= \frac{4}{\sqrt{7} \times \sqrt{4}} = \frac{4}{2 \cdot 64 \times 2} = 0.756$$

Mean $\frac{3+6}{2} = 4.5$ Similarly;

$$\text{Pearson}(U_2, U_3) = \frac{(7-4.8)(3-2) + (4-4.8)(1-2)}{(1-2) + (3-4.8)(1-2)}$$

$$= \frac{\sqrt{2.2^2 + (-0.8)^2} \times \sqrt{1^2 + (-1)^2}}{(-1.8)^2 + (-1)^2}$$

$$= \frac{4.8}{\sqrt{8.72} \times \sqrt{3}} = 0.938$$

$$\text{Pearson}(U_4, U_3) = \frac{(2-2.5)(3-2) + (2-2.5)}{(3-2) + (3-2.5)(1-2) + (3-2.5)(1-2)}$$

$$= \frac{\sqrt{(-0.5)^2 + (-0.5)^2} \times \sqrt{1^2 + 1^2 + (0.5)^2 + (0.5)^2}}{(-1)^2 + (-1)^2}$$

$$= \frac{-2}{\sqrt{1} \times \sqrt{4}} = -\frac{2}{2} = -1$$

$$\text{Pearson}(U_5, U_3) = \frac{(1-2)(3-2) + (2-2)(1-2)}{(3-2)(1-2)}$$

$$= \frac{\sqrt{(-1)^2 + 1^2} \times \sqrt{1^2 + (-1)^2}}{(-1)^2}$$

$$= \frac{-2}{\sqrt{2} \times \sqrt{2}} = -0.816$$

Hence, the top-2 closest users to user U3 are users U1 and U2 according to Pearson Correlation Coefficient.

By using the Pearson - weighted average of the raw ratings

Rank (U3, I1)

$$= \frac{(7 * 0.756 + 6 * 0.938)}{(0.756 + 0.938)} = 6.446 \Delta$$

Rank (U3, I6)

$$= \frac{(4 * 0.756 + 4 * 0.938)}{(0.756 + 0.938)} = 4 \Delta$$

Cosine Similarity;

X

$$6x3 + 7x3 + 4x1 + 5x1$$

Cosine (U1, U3) =

$$= \frac{48}{\sqrt{6^2 + 7^2 + 4^2 + 5^2} \times \sqrt{3^2 + 3^2 + 1^2}} = 0.9581$$

~~Pearson~~ (Item based)

Example 2;

Column-wise

Suppose the ratings of five users U_1, U_2, U_3, U_4 and U_5 rating is drawn from the range $\{1, \dots, 7\}$. Find the predictions of Item I_1 for user U_3 on the basis of the ratings in the table below. Use Pearson Correlation Coefficient to find similarity between users and consider top-2 closest items.

User/Item	I_1	I_2	I_3	I_4	I_5	I_6
U_1	7	6	7	4	5	4
U_2	6	7	?	4	3	4
U_3	?	3	3	1	1	?
U_4	2	2	2	3	3	4
U_5	1	?	1	2	3	3
Col ⁿ , mean	3.75	4.5	3.25	2.8	3	3.75

Solⁿ The first step is to compute the similarity between item I_1 and all the other items.

Pearson (I_1, I_2) =

$$\frac{(7-3.75)(6-4.5) + (6-3.75)(7-4.5) +}{(1-3.75)(2-4.5)}$$

$$= \frac{\sqrt{3.25^2 + 2.25^2} + \sqrt{1.5^2 + 2.5^2} + (-2.5)^2}{(-2.75)^2}$$

$$= \frac{17.375}{18.493} = 0.94$$

Pearson (I_1, I_3) =

$$= \frac{(7-3.75)(7-3.25) + (1-3.75)(2-3.25) + (1-3.75)(1-3.25)}{(-2.75)^2}$$

$$= \frac{\sqrt{3.25^2 + (-2.75)^2} + \sqrt{3.75^2 + (-1.25)^2} + (-2.25)^2}{(-2.75)^2}$$

$$= \frac{21.8125}{21.052} = 0.946$$

Pearson (I_1, I_4) =

* Content-based Recommender Systems:

Example: Find the predictions of Movie M₄ for user U₁ on the basis of the information given in the table below.

Movie/User	U1	U2	U3	U4	U5	U6	Genre	
M ₁	1	1	?	1	1	?	Romance	
M ₂	-	1	1	1	1	1	?	Thriller
M ₃	1	1	1	1	1	1	Action	
M ₄	?	1	1	1	1	1	Romance	
M ₅	1	1	1	1	1	?	Crime	
M ₆	1	1	1	1	1	1	Crime	
watched	5	6	5	6	6	3		
Item profile generation:								
movie/genre	Romance	Thriller	Action	Crime				
M ₁	1	0	0	0				
M ₂	0	1	0	0				
M ₃	0	0	1	0				
M ₄	1	0	0	0				
M ₅	0	0	0	1				
M ₆	0	0	0	1				

User profile generation:

User/Genre	Romance	Thriller	Action	Crime
U1	$\frac{1}{5} = 0.2$	0.2	0.2	$\frac{2}{5} = 0.4$
U2	0.33	0.17	0.17	0.3
U3	0.2	0.2	0.2	0.4
U4	0.33	0.17	-0.17	0.33
U5	0.33	0.17	0.17	0.33
U6	0.33	0	0.33	0.3

$$\therefore \text{Sim}(m_4, U_1) \quad x = [1, 0, 0, 0]$$

$$= \text{Cosine}(m_4, U_1) \quad x = [0.2, 0.2, 0.2, 0.4]$$

$$\theta = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}}$$

$$1 \times 0.2 + 0 \times 0.2 + 0 \times 0.2 + 0 \times 0.4$$

$$= \frac{\sqrt{1^2}}{\sqrt{0.2^2 + 0.2^2 + 0.2^2 + 0.4^2}}$$

$$= \frac{0.2}{\sqrt{0.2^2 + 0.2^2 + 0.2^2 + 0.4^2}} = 0.38 \quad \Delta$$

* Mining Frequent Patterns:

1. Apriori
2. FP growth

$$\text{Support} = \frac{\text{Frequency}}{\text{Transaction}} \times 100\%$$

1. Apriori Algorithm Confidence = $\frac{S(A \cup B)}{S(A)}$

Eg:	TID	Items
	100	1 3 4
	200	2 3 5
	300	1 2 3 5
	400	2 5

$$\text{Min Support} = 50\%$$

$$\text{Threshold Confidence} = 70\%$$

Itemset	Support (Count) / Frequency
1	2/4 \rightarrow 50%
2	3/4 \rightarrow 75%
3	3/4 \rightarrow 75%
4	1/4 \rightarrow 25% X
5	3/4 \rightarrow 75%

$$\text{Itemset} \rightarrow 1, 2, 3, 5$$

Itemset	Support
{1, 2}	$\frac{1}{4} \rightarrow 25\%$ X
{1, 3}	$\frac{2}{4} \rightarrow 50\%$
{1, 5}	$\frac{1}{4} \rightarrow 25\%$ X
{2, 3}	$\frac{2}{4} \rightarrow 50\%$
{2, 5}	$\frac{3}{4} \rightarrow 75\%$
{3, 5}	$\frac{2}{4} \rightarrow 50\%$

Itemset	Support
{1, 3, 5}	$\frac{1}{4} \rightarrow 25\%$ X
{2, 3, 5}	$\frac{2}{4} \rightarrow 50\%$
{1, 2, 3}	$\frac{1}{4} \rightarrow 25\%$ X

Rules	Support	Confidence
(2 3) \rightarrow 5	2	$2/2 = 100\%$
(3 5) \rightarrow 2	2	$2/2 = 100\%$
(2 5) \rightarrow 3	2	$2/3 = 66\%$
2 \rightarrow (2 5)	2	$2/3 = 66\%$
5 \rightarrow (2 3)	2	$2/3 = 66\%$
3 \rightarrow (2 5)	2	$2/3 = 66\%$

E.g. Confidence = $S(A \cup B) / S(A)$

$$(2 | 3) \rightarrow 5 = S(2 | 3 \cup 5) / S(2 | 3)$$

$$= 2/2 = 100\%$$

The association rules will be;

$$(2 \wedge 3) \rightarrow 5 \quad \& \quad (3 \wedge 5) \rightarrow 2 \quad \Delta$$

2. FP Growth Algorithm:

~~K~~ Example)

TID Items

1	A, B, C, D, E, F
2	B, C, D, E, F, G
3	A, D, E, H
4	A, D, F, I, J
5	B, D, E, K

$$\text{min. support} = 60\%$$

$$\text{i.e. Support Count} = 60 / 100 * 5 = 3$$

$$\text{Confidence} = 80\%$$

Step 1: 1-item set

Item	Count	Item	Count
A	3	H	1
B	3	I	1
C	2	J	1
D	5	K	1
E	4		
F	3		
G	1		

Step: 2 Item frequent set

Item	Count
A	3
B	3
D	5
E	4
F	3

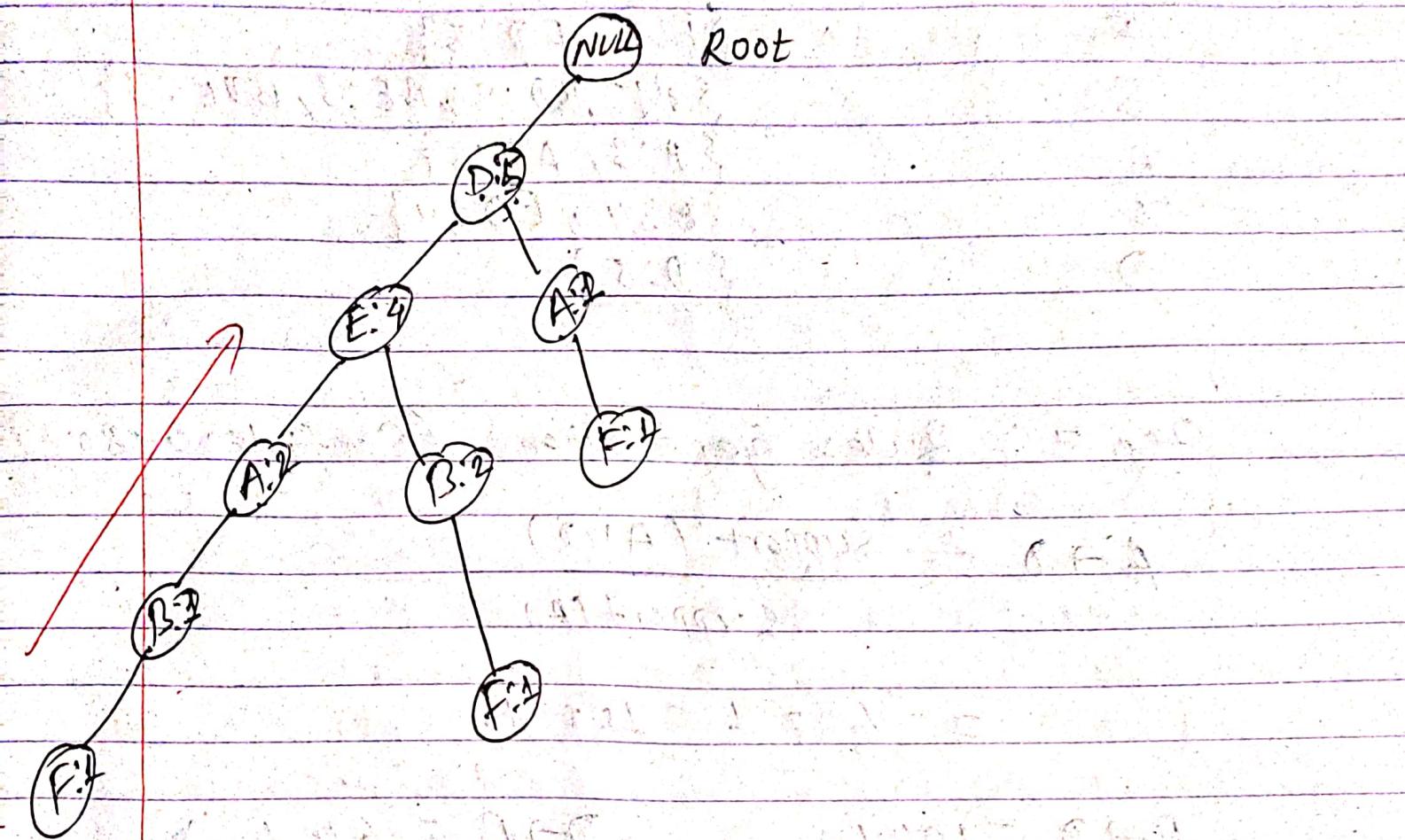
Step: 3 Recording 1-item frequent set
In descending order w.r.t. count value.

Item	Count
D	5
E	4
A	3
B	3
F	3

Step: 4 Recording original itemset in
descending order w.r.t. count value.

TID	Items
1	D, E, A, B, F
2	D, E, B, F
3	D, E, A
4	D, A, F
5	D, E, B

Step 5: Construction of FP Tree



Step 6: Generation of Frequent Patterns w.r.t Conditional pattern base.

<u>Item</u>	Conditional Pattern Base	Conditional FP Tree (Count)
-------------	--------------------------	-----------------------------

F $\{\{D, E, A, B:1\}, \{D, E, B:1\}, \{D, A:1\}\}$ $\{D:3, E:2, A:2, B:2\}$

B $\{\{D, E, A:1\}, \{D, E:2\}\}$ $\{D:3, E:3, A:1\}$

A $\{\{D, E:2\}, \{D:1\}\}$ $\{D:3, E:2\}$

E $\{D:4\}$ $\{D:4\}$

D $-$

Items

Frequent Patterns generated

F

{F: 3, FD: 3}

B

{B: 3, BD: 3, BE: 3, BDE: 3}

A

{A: 3, AD: 3}

E

{E: 4, ED: 4}

D

{D: 5}

Step 7 : Rules generation (Confidence ≥ 80%)

$$A \rightarrow D = \frac{\text{Support}(A \cup D)}{\text{Support}(A)}$$

$$= \frac{3}{3} = 1 = 100\%$$

$$B \rightarrow D = 100\%$$

$$D \rightarrow B = 60\% \times$$

$$B \rightarrow E = 100\%$$

$$E \rightarrow B = 75\% \times$$

$$D \rightarrow E = 80\%$$

$$E \rightarrow D = 100\%$$

$$D \rightarrow F = 60\% \times$$

$$F \rightarrow D = 100\%$$

$$B \rightarrow DE = 100\%$$

$$DE \rightarrow B = 100\%$$

$$D \rightarrow BE = 60\% \times$$

$$BE \rightarrow D = 100\%$$

$$E \rightarrow BD = 75\% \times$$

$$BD \rightarrow E = 100\%$$

$$E.D = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

* k-means

Q. Divide the data points $\{(1,1), (2,1), (4,3), (5,4)\}$ into two clusters.

Sol: Let $P_1 = (1,1)$, $P_2 = (2,1)$
 $P_3 = (4,3)$, $P_4 = (5,4)$

Let $C_1 = (1,1)$ and $C_2 = (2,1)$ are two initial cluster centers.

Iteration 1 Distance betw. cluster centers and each data points are;

$$d(C_1, P_1) = 0 \checkmark$$

$$d(C_2, P_1) = 1$$

$$d(C_1, P_2) = 1$$

$$d(C_2, P_2) = 0 \checkmark$$

$$d(C_1, P_3) = 3.6$$

$$d(C_2, P_3) = 2.88 \checkmark$$

$$d(C_1, P_4) = 5$$

$$d(C_2, P_4) = 4.24 \checkmark$$

Thus, after first iteration

cluster 1 = $\{P_1\}$, cluster 2 = $\{P_2, P_3, P_4\}$

Iteration 2

New cluster centers are: $C_1 (1,1)$ and $C_2 (11/3, 8/3)$.

Now, calculate distance betw. new cluster centers and each data point.

$$\begin{array}{ll}
 d(C_1, P_1) = 0 & d(C_2, P_1) = 3.14 \\
 d(C_1, P_2) = 1 & d(C_2, P_2) = 2.35 \\
 d(C_1, P_3) = 3.6 & d(C_2, P_3) = 0.46 \\
 d(C_1, P_4) = 5 & d(C_2, P_4) = 1.88
 \end{array}$$

Thus, after second iteration;

$$\text{Cluster 1} = \{P_1, P_2\}, \text{ Cluster 2} = \{P_3, P_4\}.$$

Iteration 3 New cluster Centers are:
 $C_1(3/2, 1)$ and $C_2(9/2, 7/2)$

$$\begin{array}{ll}
 d(C_1, P_1) = 0.5 & d(C_2, P_1) = 5.54 \\
 d(C_1, P_2) = 0.5 & d(C_2, P_2) = 4.33 \\
 d(C_1, P_3) = 4.5 & d(C_2, P_3) = 0.70 \\
 d(C_1, P_4) = 5.5 & d(C_2, P_4) = 0.70
 \end{array}$$

Thus, after third iteration,

$$\begin{aligned}
 \text{Cluster 1} &= \{P_1, P_2\} \\
 \text{Cluster 2} &= \{P_3, P_4\}
 \end{aligned}$$

* Robust clustering algorithm for categorical attributes (ROCK):

$$\text{Sim}(P_i, P_j) = \frac{|P_i \cap P_j|}{|P_i \cup P_j|}, \quad \theta = 0.3 \text{ (Threshold)}$$

Number of clusters = 2

$$g(P_i, P_j) = \frac{\text{Link}[P_i, P_j]}{(n+m)} = \frac{1+2f(\theta)}{n} - \frac{1+2f(\theta)}{m}$$

$$f(\theta) = \frac{1-\theta}{1+\theta}$$

Example;

$$P_1 = \{\text{judgement, faith, prayer, fair}\}$$

$$P_2 = \{\text{fasting, faith, prayer}\}$$

$$P_3 = \{\text{fair, fasting, faith}\}$$

$$P_4 = \{\text{fasting, prayer, pilgrimage}\}$$

Step 1: Similarity Table: (0-1)

	P_1	P_2	P_3	P_4
P_1	1	$\frac{2}{5} = 0.4$	$\frac{2}{5} = 0.4$	$\frac{4}{6} = 0.67$
P_2		1	$\frac{2}{4} = 0.5$	$\frac{2}{4} = 0.5$
P_3			1	$\frac{1}{5} = 0.2$
P_4				1

Taking Threshold, $\theta = 0.3$

Now,

Step 2: Adjacency Table ($\theta = 0.3$)

	P_1	P_2	P_3	P_4
P_1	1	1	1	0
P_2	1	1	1	1
P_3	1	1	1	0
P_4	0	1	0	1

(Consist of 0's & 1's) $\text{Similarity} > \theta = 0.3$
 $\Rightarrow 1$ otherwise 0.

Step 3: Link Table ($A \times A$)
 (Row * column) multiplication

	P_1	P_2	P_3	P_4	
P_1	-	3	3	1	x
P_2	-	3	2		
P_3		-	1		
P_4				-	

Pair

Goodness Measure

$\{P_1, P_2\}$	1.3825	Search for highest values.
$\{P_1, P_3\}$	1.3825	
$\{P_1, P_4\}$	0.46	
$\{P_2, P_3\}$	1.3825	
$\{P_2, P_4\}$	0.92	
$\{P_3, P_4\}$	0.46	

$$g(P_1, P_2) = \frac{\text{Link}(P_1, P_2)}{(n+m) 1 + 2f(0)} = \frac{1 + 2f(0)}{n + m} = 1 + 2f(0)$$

$$\text{Assuming } f(0) = \frac{1 - 0}{1 + 0} = \frac{1 - 0.3}{1 + 0.3} = 0.538$$

~~Ans~~,

$$= \frac{3}{(1+1)^{1+2 \times 0.5^3} - 1 - 1}$$

$$= \frac{3}{2^{2.06} - 2} = 1.3825$$

For (P_1, P_3) ;

$$= \frac{3}{(1+1)^{1+2 \times 0.5^3} - 1 - 1} = 1.3825$$

For (P_1, P_4) ;

$$= \frac{1}{2^{2.06} - 2} = 0.46$$

For (P_2, P_3) ;

$$= \frac{3}{2^{2.06} - 2} = 1.3825$$

For (P_2, P_4) ;

$$= \frac{2}{2^{2.06} - 2} = 0.921$$

For (P_3, P_4) ;

$$= \frac{1}{2^{2.06} - 2} = 0.46$$

The pairs (P_1, P_2) , (P_1, P_3) and (P_2, P_3) have highest goodness measure among all other pairs. So, next pair will be;

<u>Link Table</u>	Pair $\{P_1, P_2\}$	P_3	P_4	<u>Link Table</u>
$\{P_2, P_3\}$	—	$3+3=6$	$1+2=3$	
P_3	—	—	1	
P_4	—	—	—	

<u>Pair</u>	<u>Goodness Measure</u>
$\{P_1, P_2\}, P_3$	• 1.35 ✕
$\{P_1, P_2\}, P_4$	0.67
$\{P_2, P_3\}$	0.46

For $\{P_1, P_2\}, P_3$;

$$= \frac{6}{(2+1)^{2.06} - 2 - 1} = 1.35$$

For $\{P_1, P_2\}, P_4$

$$= \frac{3}{(2+1)^{2.06} - 2^{2.06} - 1} = 0.67$$

For $\{P_3, P_4\}$,

$$= \frac{1}{(1+1)^{2.06} - 1^{2.06} - 1} = 0.46$$

$\{P_1, P_2, P_3\}$ ~~P₄~~

$\{P_1, P_2, P_3\}$ — $3+1=4$

$\{P_4\}$ —

Pair Goodness measure

$\{P_1, P_2, P_3\}, P_4$ 0.59

$$= \frac{4}{(3+1)^{2.06} - 3^{2.06} - 1}$$

$$= 0.59$$

There is only two clusters. So, stop + process.

* DBSCAN:

DBSCAN is a density-based clustering algorithm that groups together points that are closely packed together while marking points in low-density regions as outliers.

It requires two parameters: $\text{epsilon}(\epsilon)$, which defines the radius of the neighbourhood around a point, and minpts , which specifies the minimum number of points required to form a dense region.

→ It stands for Density-Based spatial Clustering of Applications with Noise.

E.g. $P_1: (3, 7)$, $P_2: (4, 6)$; $P_3: (5, 5)$, $P_4: (5, 4)$,
 $P_5: (7, 3)$

	P_1	P_2	P_3	P_4	P_5
P_1	0				
P_2	1.41	0			
P_3	2.83	1.41	0		
P_4	4.24	2.83	1.41	0	
P_5	5.66	4.24	2.83	1.41	0

Given,

$$\text{minPts} = 3$$

$$\text{epsilon}(\varepsilon) = 1.9$$

Step: 1 Find distance less than or equal to $\text{epsilon}(\varepsilon)$; i.e., 1.9.
(Search horizontally & vertically).

For:

$$P_1 : P_2, \cancel{P_3} \quad P_3 : P_2, P_9$$

$$P_2 : P_1, P_3 \quad P_4 : P_3, P_5$$

& so on. : P_5 : P_4.

Point	Status
P ₁	Noise
P ₂	Core
P ₃	Core
P ₄	Core
P ₅	Noise

P₁ should be a part of anyone core data points.

3 data centers
(3 core)

* Bayesian Classification:

Bayesian Classification is based on Bayes' theorem. It is also called Naive Bayes classification or Naive Bayesian classification. Bayes' Theorem is given by:

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

Bayes' theorem is useful in that it provides a way of calculating the posterior probability, $P(H|X)$, from $P(H)$, $P(X|H)$, and $P(X)$. Here $P(X)$ and $P(H)$ are prior probability.

Let D be a database and C_1, C_2, \dots, C_m are m classes. Now above Bayes rule can be written as:

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

Let X is the set of attributes $\{x_1, x_2, x_3, \dots, x_n\}$ where attributes are independent of one another. Now, the probability $P(X|C_i)$ is given by the equation given below:

$$P(X|c_i) = \prod_{k=1}^n P(x_k|c_i)$$

$$= P(x_1|c_i) \times P(x_2|c_i) \dots \times P(x_n|c_i)$$

Example:

buys-Computer

RID	age	income	student	credit	class
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle-aged	high	no	fair	yes
4	Senior	medium	no	fair	yes
5	Senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7.	middle-aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	Senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle-aged	medium	no	excellent	yes
13	middle-aged	high	yes	fair	yes
14	senior	medium	no	excellent	no
15.	youth	medium	yes	fair	?

Predict class level of the tuple: x

= (age = youth, income = medium, student = yes, credit-rating = fair) using Bayesian classification.

SOLN: Prior probability of each class can be computed based on the training triples:

$$P(\text{buys-Computer} = \text{yes}) = 9/14 = 0.643$$

$$P(\text{buys-Computer} = \text{no}) = 5/14 = 0.357$$

Compute the following conditional probabilities:

$$\times P(\text{age} = \text{youth} | \text{buys-Computer} = \text{yes}) = 2/9 = 0.222$$

$$P(\text{age} = \text{youth} | \text{buys-Computer} = \text{no}) = 3/5 = 0.6$$

$$\times P(\text{income} = \text{medium} | \text{buys-Computer} = \text{yes}) = 4/9 = 0.444$$

$$P(\text{income} = \text{medium} | \text{buys-Computer} = \text{no}) = 2/5$$

$$\times P(\text{student} = \text{yes} | \text{buys-Computer} = \text{yes}) = 0.4$$

$$= 6/9 = 0.667$$

$$P(\text{student} = \text{yes} | \text{buys-Computer} = \text{no}) = 1/5 = 0.2$$

$$\times P(\text{Credit-rating} = \text{fair} | \text{buys-Computer} = \text{yes}) = 6/9 = 0.667$$

$$P(\text{Credit-rating} = \text{fair} | \text{buys-Computer} = \text{no}) = 2/5 = 0.4$$

Using the above probabilities, we obtain

$$P(X \mid \text{buys-computer} = \text{yes})$$

$$= P(\text{age} = \text{youth} \mid \text{buys-computer} = \text{yes}) \times \\ P(\text{income} = \text{medium} \mid \text{buys-computer} = \text{yes}) \times \\ P(\text{student} = \text{yes} \mid \text{buys-computer} = \text{yes}) \times \\ P(\text{Credit-rating} = \text{fair} \mid \text{buys-computer} = \text{yes})$$

$$= 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$$

$$P(X \mid \text{buys-computer} = \text{no})$$

$$= P(\text{age} = \text{youth} \mid \text{buys-computer} = \text{no}) \times \\ P(\text{income} = \text{medium} \mid \text{buys-computer} = \text{no}) \times \\ P(\text{student} = \text{yes} \mid \text{buys-computer} = \text{no}) \times \\ P(\text{Credit-rating} = \text{fair} \mid \text{buys-computer} = \text{no})$$

$$= 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$$

Now,

$$\underline{\underline{P(\text{buys-computer} = \text{yes} \mid X)}}$$

$$= P(X \mid \text{buys-computer} = \text{yes}) \times P(\text{buys-computer} = \text{yes}) \\ = 0.044 \times 0.643 = 0.028$$

$$P(\text{buys-computer} = \text{no} \mid X)$$

$$= P(X \mid \text{buys-computer} = \text{no}) \times P(\text{buys-computer} = \text{no})$$

ID3 Decision Tree

ID3 Stands for Iterative Dichotomiser 3. It uses top-down greedy approach to build decision tree model.

This algorithm computes information gain for each attribute and then selects the attribute with the highest information gain.

Information gain measures reduction in entropy after data transformation. It is calculated by comparing entropy of the dataset before and after transformation.

Entropy is the measure of homogeneity of the sample. Entropy or expected information of dataset D is calculated by using equation 1.

$$E(D) = - \sum_{i=1}^m p_i \log_2 p_i \quad \dots \quad (1)$$

where,

p_i is the probability of a tuple in D belonging to class C_i and is estimated using Equation 2.

$$p_i = \frac{|C_i, D|}{|D|} \quad \dots \quad (2)$$

where, $|C_i, D|$ is the number of tuples in D belonging to class C_i and $|D|$ is the number of tuples in D .

Suppose we have to partition the tuples in D on some attribute A having r distinct values. The attribute A can be used to split D into r partitions $\{D_1, D_2, \dots, D_r\}$.

Now, Entropy of attribute, A is calculated as,

$$E_A(D) = \sum_{j=1}^r \frac{|D_j|}{|D|} \times E(D_j)$$

Finally, the information gain achieved after partitioning D on attribute A is calculated as,

$$IG(A) = E(D) - E_A(D).$$

Example:

Buyer

	age	income	student	credit-rating	class
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle-aged	high	no	fair	yes
4	Senior	medium	no	fair	yes
5	Senior	low	yes	fair	yes
6	Senior	low	yes	excellent	no
7	middle-aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	Senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle-aged	medium	no	excellent	yes
13	middle-aged	high	yes	fair	yes
14	Senior	medium	no	excellent	no

Construct decision tree from above data using information gain.

- Predict class level of the tuple: $x = (\text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit-rating} = \text{fair})$.

soln:

- Expected information needed to classify a tuple in D is:

$$E(D) = -\frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right)$$

$$= 0.94$$

$$\text{Eage}(D) = \frac{5}{14} \left[-\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right] +$$

$$\frac{4}{14} \left[-\frac{4}{4} \log_2 \left(\frac{4}{4} \right) - \frac{0}{4} \log_2 \left(\frac{0}{4} \right) \right] + \frac{5}{14} \left[-\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right]$$

$$= 0.694$$

~~$$\therefore \text{IG(Age)} = E(D) - \text{Eage}(D)$$~~

$$= 0.94 - 0.694 = 0.246$$

Similarly,

$$\text{Eincome}(D) = \frac{4}{14} \left[-\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right]$$

$$+ \frac{6}{14} \left[-\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \right] + \frac{4}{14}$$

$$\left[-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right] = 0.01$$

$$\text{IG(Income)} = E(D) - E(\text{income}(D))$$

$$= 0.94 - 0.91 = 0.03$$

Again,

$$E(\text{student}(D)) = \frac{7}{14} \left[-\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} \right] +$$

$$\frac{7}{14} \left[-\frac{3}{7} \log_2 \frac{3}{7} + \frac{4}{7} \log_2 \frac{4}{7} \right]$$

$$= 0.2958 + 0.4926 = 0.788$$

$$\text{IG(Student)} = E(D) - E(\text{student}(D))$$

$$= 0.94 - 0.788 = 0.152$$

$$E(\text{credit_rat}(D)) = \frac{8}{14} \left[-\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \right]$$

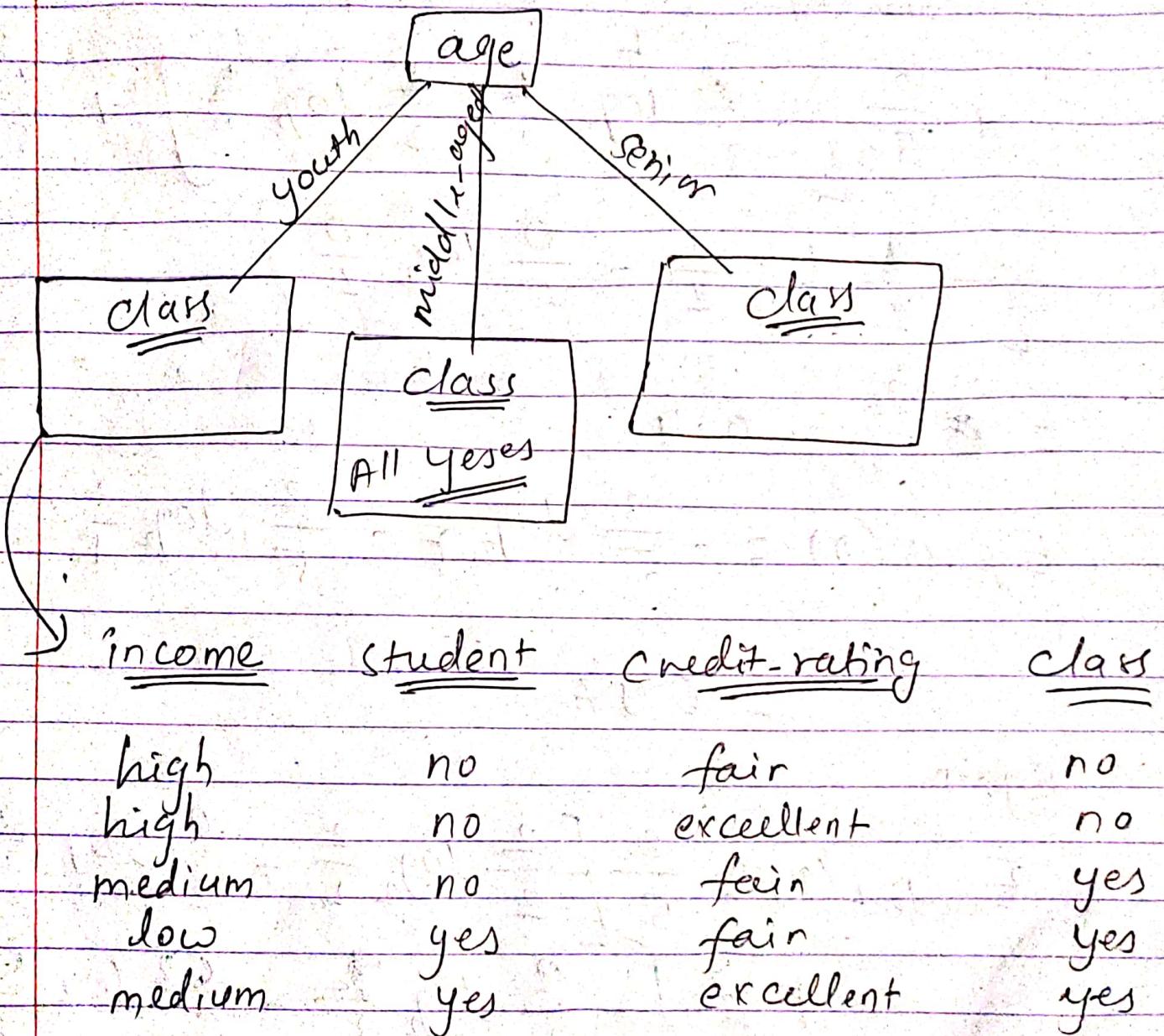
$$+ \frac{6}{14} \left[-\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \right]$$

$$= 0.4635 + 0.4285 = 0.892$$

$$\text{IG(Credit_rat)} = E(D) - E(\text{credit_rat}(D))$$

$$= 0.94 - 0.892 = 0.048$$

Because age has the highest information gain among the attributes, it is selected as the splitting attribute. The decision tree now looks like below:



Now, calculate information gain of income, student and credit-rating in case of left children (youth) & right children (Senior). And repeat the above process.

$$\left. \begin{array}{l} \log 0 = 0 \\ \log 1 = 0 \end{array} \right\}$$

$$IG(\text{income}) = E_{\text{youth}}(D) - E_{\text{income}}(D)$$

~~0.57~~

$$= \cancel{0.691} - 0.97 - 0.1 = \cancel{0.97}$$

$$E_{\text{income}}(D) = \frac{2}{5} \left[-\frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log_2 \frac{2}{2} \right]$$

$$+ \frac{2}{5} \left[-\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} \right] + \frac{1}{5}$$

$$\left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{0}{1} \log_2 \frac{0}{1} \right]$$

$$= \cancel{\frac{2}{5}} = 0.4$$

$$E_{\text{youth}}(D) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}$$

$$= 0.97$$

Again,

$$E_{\text{student}}(D) = \frac{2}{5} \left[-\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} \right]$$

$$+ \frac{3}{5} \left[-\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3} \right]$$

$$= \cancel{0.5509} \quad 0$$

$$IG(\text{student}) = E_{\text{youth}}(D) - E_{\text{student}}(D)$$

$$= 0.97 - \cancel{0.5509} = \cancel{0.4191}$$

$$E_{\text{credit-rat}}(D) = \frac{3}{5} \left[-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right] + \frac{2}{5} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right]$$

$$= 0.57097 + 0.4 = 0.950$$

Q $I_G(\text{credit-rat}) = E_{\text{youth}}(D) - E_{\text{credit-rat}}(D)$

$$= 0.97 - 0.95 = 0.02$$

Because ~~income~~ has the highest information gain.