

Chapt 5. Data Mining

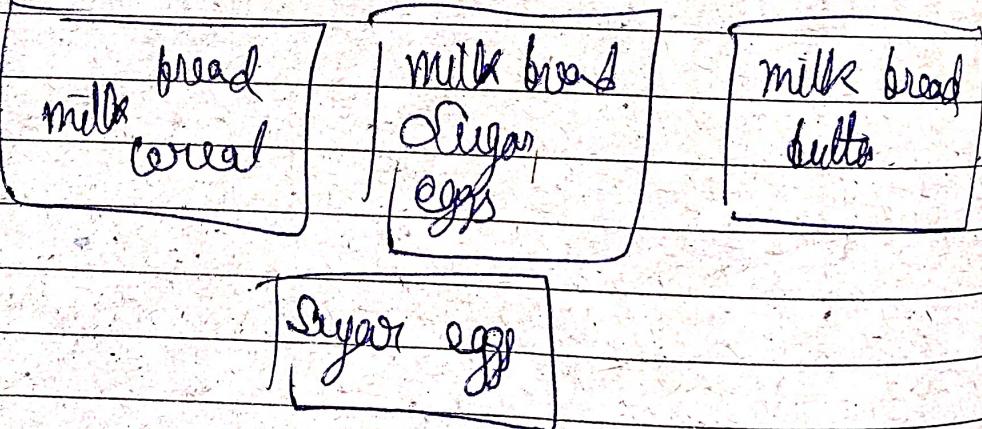
Frequent Patterns → Freq patterns are factors that appear frequently in a data set.

A typical example of frequent itemset mining is market basket analysis. This process analyzes customer buying habits by finding associations b/w the different items that customers place in their "shopping baskets".

The discovery of those association can help retailer develop marketing strategies by gaining insight into which items are frequently purchased together by customers.

Goal: To find items are frequently purchased together by customer

market analyst



market based analysis

Items that are frequently purchased together can be placed in proximity to further encourage the combined sale of such items.

Placing hardware & software at opposite ends of the store may entice customers who purchase such items to pick up other items along the way.

Association Rules

Customers who purchase Computer also tend to buy anti-virus software at the same time.

Computer \Rightarrow anti-virus Software [Support = 2%, Confidence = 60%]

* Support and Confidence are two measures of rule interestingness.

Support \rightarrow Replace the usefulness

Confidence \rightarrow Certainty of discovered rule

* A support of 2% means that 2% of all the transaction under analysis show that Computer and anti-virus software are purchased together.

* A confidence of 60% means that 60% of the customers who purchased a Computer also bought this software.

Typically, association rules are considered interesting if they satisfy both a minimum support threshold and a minimum confidence threshold.

These thresholds can be set by users or domain experts.

$$\text{Support } (A \rightarrow B) = P(A \cup B)$$

$$\text{Confidence } (A \rightarrow B) = P(B|A)$$

Rules that satisfy both minimum support threshold and a minimum confidence threshold are called strong.

The occurrence frequency of an itemset is the number of transactions that contain the itemset. This is also called frequency, Support Count, Count.

Note \rightarrow $\text{Support } (A \rightarrow B) = P(A \cup B) \rightarrow$ relative support

Occurrence frequency \rightarrow absolute support

If the relative support of an itemset I satisfies a minimum support threshold (i.e., the absolute support of item I satisfies a minimum support count threshold) then I is a frequent itemset.

$$\text{Confidence } (A \rightarrow B) = P(B|A) = \frac{\text{Support } (A \cup B)}{\text{Support } (A)}$$

$$= \frac{\text{Support Count } (A \cup B)}{\text{Support Count } (A)}$$

In general, association rule mining can be illustrated as a two step process :

Find all frequent itemsets : By definition, each of these itemset will occur at least as frequently as a predetermined minimum support count, min-sup

Generate strong association rules from the frequent itemsets : By definition, these rules must satisfy minimum support and minimum confidence.

Because the second step is much less costly than the first, overall performance of mining association rules is determined by the first step.

A major challenge in mining frequent itemsets from large data set is the fact that such mining often generate a huge number of itemsets satisfying the min-sup criteria especially when min-sup is low.

for ex. 100 buyer of itemset frequent

$${}^{100}C_1 = 100 \text{ frequent 1-itemset } \{q_1, q_2\}, \dots$$

$${}^{100}C_2 = \text{frequent 2-itemset } \{q_1, q_2\}, \{q_1, q_3\}, \dots$$

$$\text{Total} \Rightarrow {}^{100}C_1 + {}^{100}C_2 + {}^{100}C_3 + \dots {}^{100}C_{100} = 2^{100} - 1$$

$$\approx 1.27 \times 10^{30}$$

This is too huge number of item set for any computer to complete in time.

To over come this difficulty, we introduce the concept of closed frequent itemset and maximal frequent itemset

① Closed frequent itemset : \rightarrow An itemset is closed if none of its immediate super~~subset~~^{support} has same support as that of itemset.

Or

"An itemset X is closed in dataset D iff there exists no proper super-itemset Y such that Y has the same support count as X in D"

② Maximal frequent itemset : \rightarrow An itemset is maximal if none of its immediate ~~subset~~^{super} superset is frequent.

"An itemset X in dataset D if X is frequent, and there exists no super itemset Y such that $X \subset Y$ and Y is frequent in D."

item[Count]

A(3)

B(4)

C(5)

D(4)

$$\text{min-sup} = 3$$

	Item	Count
A, B	3	
A, C	3	
A, D	2	not frequent
B, C	4	
B, D	3	
C, D	4	

A(3) \rightarrow AB(3), AC(3), AD(2)

A(count) is not greater than its immediate superset.
A is not closed.

In A's Immediate superset, itemset A present with
min-sup count = 3.

A is not maximal.

Apriori Property

Employs an iterative approach known as a level-wise search, where K-itemsets are used to explore (K+1)-itemsets.

To improve the efficiency of level wise generation of frequent itemsets, an important property called the Apriori Property -

→ All nonempty subsets of a frequent itemset must also be frequent.

Antimonotonicity → if a set passes a test, all its supersets will fail the same test as well.

Min-Suffix = 2

Date: _____
Page No.: _____

C1

<u>Ex</u>	<u>TID</u>	<u>Items</u>	<u>Itemset</u>	<u>Sub Count</u>
	T100	1, 2, 5	1	6.
	T2	2, 4	2	7
	T7	2, 3	3	6
	T4	1, 2, 4	4	2
	T5	1, 3	5	2
	T6	2, 3	L1	
	T7	1, 3		
	T8	1, 2, 7, 5	Itemset	Sub Count
	T9	1, 4, 3	1	6
			2	7
			3	6
			4	2
			5	2

C2

<u>Items</u>	<u>Sub-count</u>	<u>Itemset</u>	<u>Sub-Count</u>
1, 2}	4		
1, 3}	4		
1, 4}	1	X	L2
1, 5}	2		
2, 3}	4	2, 1, 3}	4
2, 4}	2	2, 1, 3}	4
2, 5}	2	2, 1, 5}	2
3, 4}	0	X	2, 3}
3, 5}	0	X	2, 4}
4, 5}	0	X	2, 5}

C2 / L3

<u>Itemset</u>	<u>Sub-Count</u>
1, 2, 3}	2
1, 2, 5}	2

			Support	Date:	Page No.	Confidential
1, 2, 3	=	1 ∩ 2 → 3	2	2/4	50%	
		1 ∩ 3 → 2	2	2/4	50%	
		2 ∩ 3 → 1	2	2/4	50%	
1, 2, 5	=	1 ∩ 2 → 5	2	2/4	50%	
		1 ∩ 5 → 2	2	2/2	100%	
		2 ∩ 5 → 1	2	2/2	100%	
		3 → 1 ∩ 2	2	2/6	33.3	
$I(A \rightarrow B) =$	<u>Support</u>	2 → 1 ∩ 3	2	2/7	28.5	
	<u>AP(A)</u>	1 → 2 ∩ 3	2	2/6	33.3	
	\Rightarrow	5 → 1 ∩ 2	2	2/2	100	
		2 → 1 ∩ 5	2	2/7	28.5	
		1 → 1 ∩ 5	2	2/6	33.3	

Disadvantages

Requires Many Database Scans

Assumes Transactions Data has in Memory Residency

Improving the Efficiency of Apriori

I Hash Based Technique →

ID	Key	Itemset	Support Count
P1	1, 2, 5	1	6
P2	2, 4	2	7
P3	2, 3	3	6
P4	1, 2, 4	4	2
P5	1, 3	5	2
P6	2, 3		
P7	1, 3		
P8	1, 2, 3, 5		
P9	1, 4, 3		

$h(x|y) = ((\text{order}(x) \times 10 + (\text{ord}(y))) \mod 7)$

skew | Count | Hash f⁴

1, 2	4	$(1 \times 10 + 2) \% 7 = 5$
1, 3	4	$(1 \times 10 + 3) \% 7 = 6$
1, 4	1	$(1 \times 10 + 4) \% 7 = 0$
1, 5	2	$(1 \times 10 + 5) \% 7 = 1$
2, 3	4	$(2 \times 10 + 3) \% 7 = 2$
2, 4	2	$(2 \times 10 + 4) \% 7 = 3$
2, 5	2	$(2 \times 10 + 5) \% 7 = 4$
3, 4	0	$(3 \times 10 + 4) \% 7 =$
3, 5	1	$(3 \times 10 + 5) \% 7 = 0$
4, 5	0	

Hash Tab H₂

Bucket address	0	1	2	3	4	5	6
Bucket Count	1+1	2	4	2	2	9	4
Bucket Content	{1, 2}	{1, 3}	{2, 3}	{1, 2, 3}	{1, 2, 3}	{1, 2, 3}	{1, 2, 3}
	{3, 4, 5}	{1, 5}	{2, 4, 5}	{2, 4, 5}	{2, 4, 5}	{1, 2, 3}	{1, 2, 3}

Transaction Reducer \rightarrow A transaction that does not contain any frequent k items, cannot contain any frequent $(k+1)$ items. Therefore, such a transaction can be removed from further consideration because subsequent database scan for j items, where $j > k$, will not need to consider.

Ex

T1	1, 2, 5
T2	2, 3, 4
T3	3, 4
T4	1, 2, 3, 4

$$\min_sup = ?$$

	1	2	3	4	5		1	2	3	4	
T1	1	1	0	0	1		T1	1	1	0	0
T2	0	1	1	1	0		T2	0	1	1	1
T3	0	0	1	1	0		T3	0	0	1	1
T4	1	1	1	1	0		T4	1	1	1	1

(*)

	2{1,2,3}	2{1,3}	2{1,4,3}	2{2,3}	2{2,4,3}	2{3,4,3}
T1	1	0	0	0	0	0
T2	0	0	0	1	1	0
T3	0	0	0	0	0	1
T4	1	1	1	1	1	1

	2{1,2,3}	2{2,3}	2{2,4,3}	2{3,4,3}
T2	0	1	1	0
T4	1	1	1	1

↓

	2{1,2,3}	2{1,2,4,3}	2{2,3,4,3}
T2	0	0	1
T4	1	1	1

↓

	2{2,3,4,3}
T4	1

↓

Partitioning

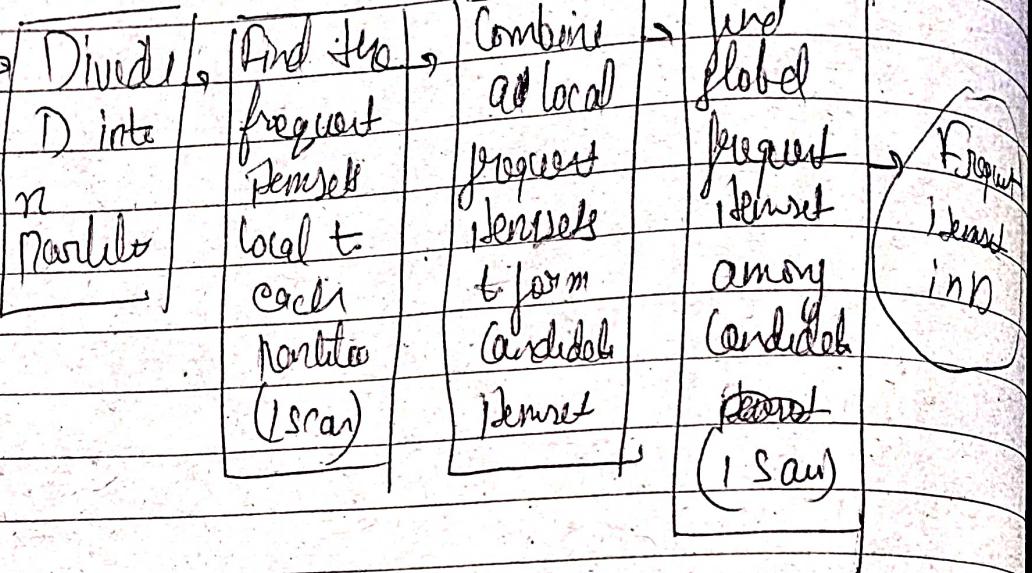
Any itemset that is Potentially Frequent in DB must be frequent in at least one of the partitions of DB (2 DB Scan)

Phase I

Date : _____
Page No. : _____

Phase II

Transaction
in D



Ep 1 2 3 4 5

T1	1	0	0	0	1
T2	0	1	0	1	0
T3	0	0	0	1	1
T4	0	1	1	0	0
T5	0	0	0	0	1
T6	0	1	1	1	0

Data base is divided into
3 Partitions

Each having 20%
Min-Sup

Trans	1,2,3	First Scan. Sup = 20% minSup = 1	Second Scan Sup = 20% minSup = 2	Output
T1	1,5	1 → 1, 2 → 1	1 → 1, 2 → 3	
T2	2,4	4 → 1, 5 → 1, {1,5} → 1, {2,4} → 1	3 → 2, 4 → 3	2 → 3
T3	4,5	4 → 1, 5 → 1, 2 → 1, 3 → 1	5 → 3,	3 → 2
T4	2,3	{4,5} → 1, {2,3} → 1	{1,5} → 1	4 → 3
T5	5	5 → 1, 2 → 1, 3 → 1, 4 → 1	{4,5} → 1	5 → 3
T6	2,3,4	{2,3} → 1, {2,4} → 1	{2,4} → 2	{2,4} → 2
			{4,5} → 1	{2,3} → 2
			{1,3} → 2	
			{1,4} → 1	
			{1,3,4} → 1	

Dynamic Itemset Counting

It is an algorithm which reduces the number of passes made over the data while keeping the number of itemsets which are counted in any pass relatively low.

This technique can add new candidate itemset at any marked start point of the database during the scanning of the DB.

Sampling

The basic idea is to pick up a random sample S of the given data D , and then search for frequent itemsets in S instead of D .

It may be possible to lose a global frequent itemset.

This can be reduced by lowering the Min. Support.

There is Trade off some degree. \Rightarrow Accuracy vs. computation efficiency.

FPGrowth algo

Divide & Conquer strategy

Ex min Supp/len = 30%

Tid	Item
1	E, A, D, B
2	D, A, C, E, B
3	C, A, B, E
4	B, A, D
5	D
6	D, D
7	A, D, E
8	B, C

$$\min \{ \text{Pr}^n \} = \frac{8 \times 30}{100} = 2.4 \approx 3$$

Item	Sup Count	Priority
A	5	3
B	6	1 st
C	3	5
D	6	2
E	4	4

Tid	Item	ordered Items	
1	E, A, D, B	B, D, A, E	B: 1, 2, 3, 4, 5
2	D, A, C, E, B	B, D, A, E, C	D: 1, 2, 3, 4
3	C, A, B, E	B, A, E, C	A: 1, 2, 3
4	B, A, D	B, D, A	E: 1, 2
5	D	D	C: 1
6	D, B	B, D	A: 1
7	A, D, E	D, A, E	E: 1
8	B, C	B, C	C: 1

A: 1
 E: 1
 C: 1

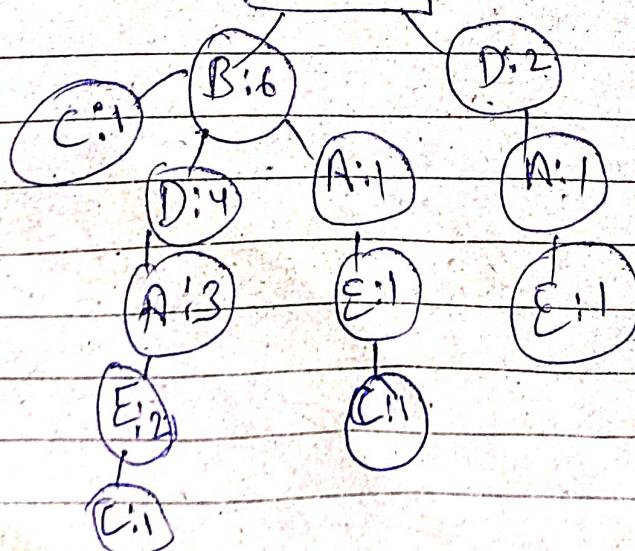
D: 1, 2

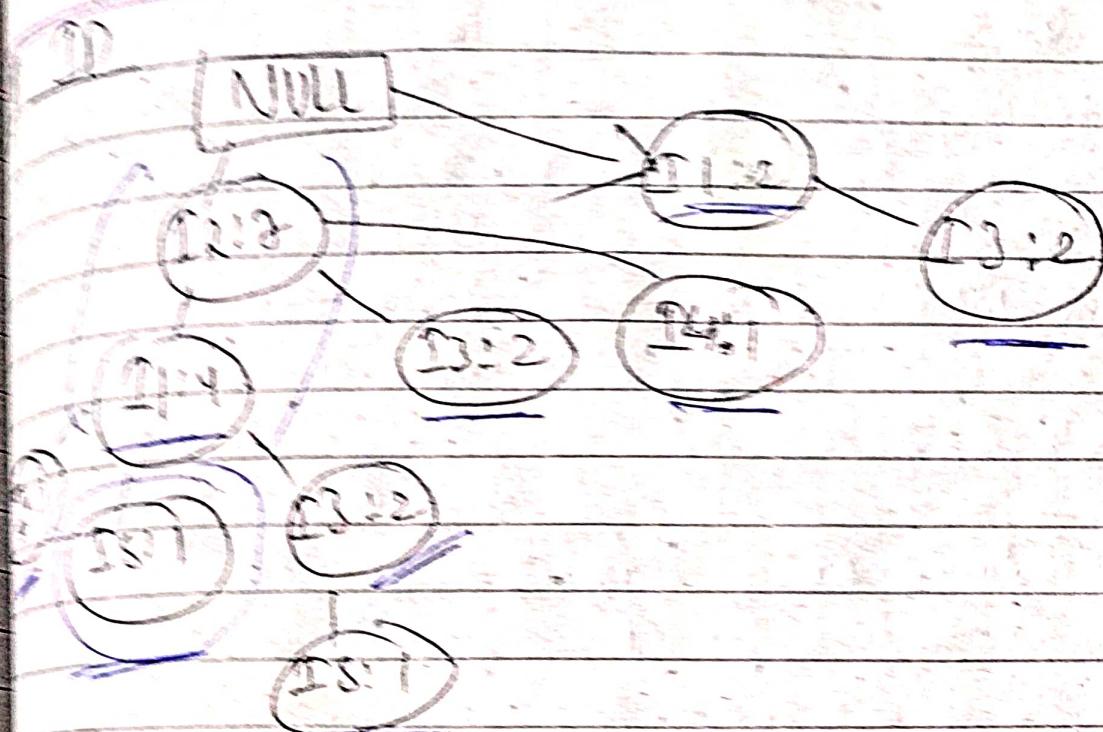
A: 1

E: 1

C: 1

NULL





1. Conventional Pattern Base

15	$\{12, 21:1\} \{12, D, B:13\}$	form Qu in Jump French
16	$\{2, 21:2\} \{22:2\}$	
17	$\{1, 21:2\} \{P2:2\} \{21:23\} \{22:3\} \{23\} \{1D\}$	
21	$\{22:43\}$	

Conventional A True

18	$\{12, 21:23\} \{22:2\}$	Free Pattern Generated
22	$\{22:2\}$	
23	$\{22:23\} \{21:28\}$	
24	$\{22:43\}$	

Free Pattern Generated

19	$\{22, 25:23\} \{21, 21:2\} \{23, 21, 25:23\}$
20	$\{22, 24:2\}$
21	$\{22, 23:23\} \{21, 23:43\} ; \{22, 21, 23:23\}$
22	$\{22, 21:43\}$

TID | Itemset

- 1 $\{a, d, e\}$
- 2 $\{a, d, b, c\}$
- 3 $\{a, d, b\}$
- 4 $\{a, e\}$
- 5 $\{b, c\}$
- 6 $\{a, d, b, e, c\}$

$$min - supp = 2$$

Item

Support Count

Priority

a

5

1

b

4

2

c

3

4

d

4

3

e

3

5

TID

Itemset

Ordered itemset

a:1, 2, 3, 4, 5

d:1

1

$\{a, d, e\}$

$\{a, d, e\}$

e:1

2

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

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

e:1

3

$\{a, d, b\}$

$\{a, b, d\}$

b:1, 2, 3

4

$\{a, e\}$

$\{a, e\}$

a:1, 2, 3

5

$\{b, c\}$

$\{b, c\}$

c:1, 2

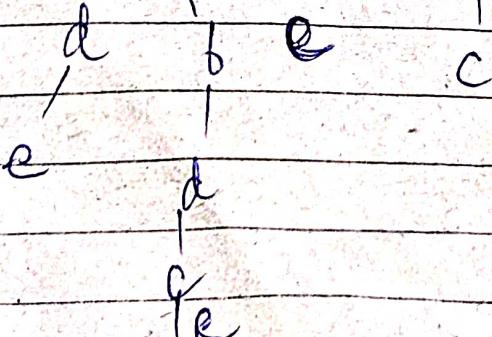
6

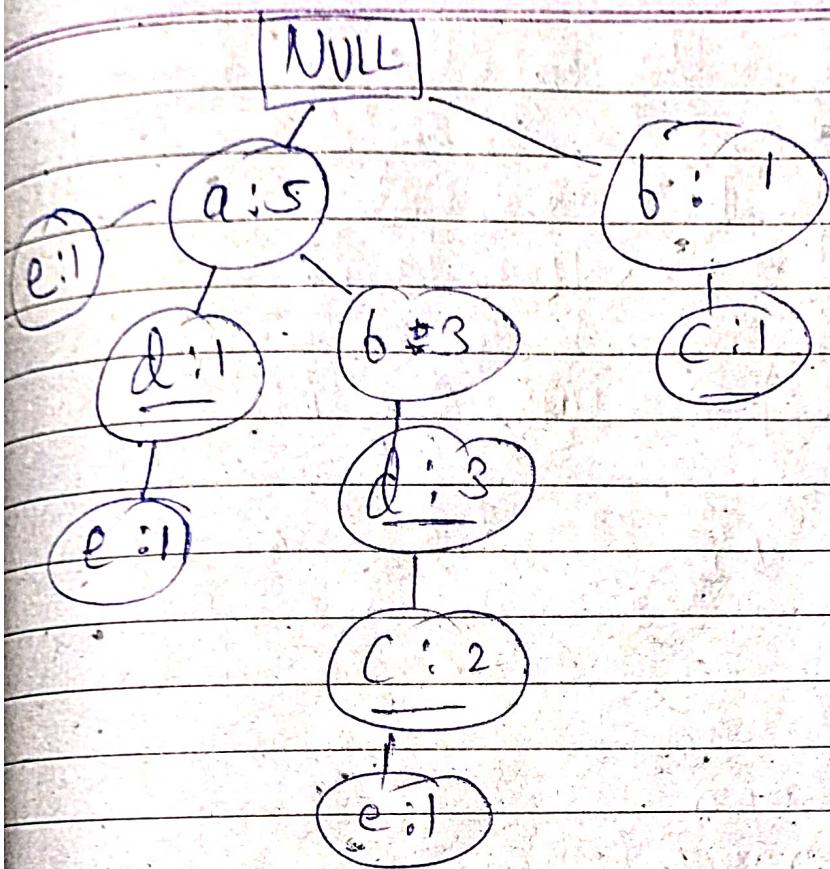
$\{a, d, b, e, c\}$

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

e:1

NULL





Item | Conditional Pattern base

e | {a, b, d, c: 1} ∪ {a, d: 1}, {a: 1}

all are in same branch

c | {a, b, d: 2} ∪ {b: 1}

all are in same branch

d | {a: 1} ∪ {a, b, c: 3}

b | {a: 3}

Item | Conditional FP tree

e | {a: 3, e: 2}

c | {a, b, d: 2}

d | {a: 4, b: 3, e: 3}

b | {a: 3}

Freq. Patterns Generated

{a, e: 3} ∪ {d, e: 2} ∪ {a, d, e: 2}

{a, c: 2} ∪ {b, c: 2} ∪ {d, c: 2} ∪ {a, b, c: 2}

{a, d, e: 2} ∪ {b, d, c: 2} ∪ {a, b, d, c: 2}

{a, d: 4} ∪ {b, d: 3} ∪ {a, b, d: 3}

{a, b: 3}

Horizontal format

$T_1 : \{I_1, I_2, I_3\}$

$T_2 : \{I_2, I_5\}$

Vertical Format

$I_1 : \{T_{100}, T_{200}\}$

$I_2 : \{T_{100}, T_{300}\}$

buys (X , "Computer Games") \Rightarrow buys (X , "Videos")

[Support = 40%, Confidence = 60%]

Total transaction = 10000

Computer games = 6000

Videos = 2500

Computer game (A) Video (B) = 4000

$$\text{Support}(A \rightarrow B) = \frac{P(A \cup B)}{P(A)} = \frac{4000}{10000} = 40\%$$

$$\text{Confidence}(A \rightarrow B) = \frac{P(A \cup B)}{P(A)} = \frac{4000}{6000} = \frac{2}{3} = 66.7\%$$

left

$$\text{lift}(A, B) = \frac{P(A \cup B)}{P(A) P(B)} = \frac{\text{Sup}(A \cup B)}{\text{Sup}(A) \text{Sup}(B)}$$

$\text{lift}(A, B) < 1 \rightarrow$ negatively correlated

$\text{lift}(A, B) > 1 \rightarrow$ positively correlated

$\text{lift}(A, B) = 1 \rightarrow$ Independent

$$\text{lift}(A, B) = \frac{P(B|A)}{P(B)} = \frac{\text{Confidence}(A \rightarrow B)}{\text{Sup}(B)}$$

χ^2 Correlation

	game	game	
Video	4000 (4500)	2500 (3000)	7500
Video	2000 (1500)	500 (400)	2500
	6000	4000	10000

$$\text{Expected } E_{ij} = \frac{\sum \text{row}_i}{\text{Total}} \times \frac{\sum \text{col}_j}{\text{Total}}$$

$$E_{11} = \frac{2500 \times 6000}{10000} = 1500$$

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}} = \frac{(4500 - 4000)^2}{4000} + \frac{(3000 - 2500)^2}{2500} + \frac{(1500 - 1000)^2}{1000} + \frac{(4000 - 3500)^2}{3500}$$

$$= 555.6$$

Because the χ^2 value is greater than 1, and the observed value of the slot (game, video) = 4000 which is less than expected value of 4500.

So buying game and buying video are negatively correlated.

Pattern Evaluation Measures

1) all Confidence (A, B)

$$= \frac{\text{Sup}(A \cup B)}{\max\{\text{Sup}(A), \text{Sup}(B)\}} = \min\{P(A|B), P(B|A)\}$$

2) max-conf (A, B) =

$$\max\{P(A|B), P(B|A)\}$$

~~$$= \max\left\{ \frac{\text{Sup}(A \cup B)}{\text{Sup}(A)}, \frac{\text{Sup}(A \cup B)}{\text{Sup}(B)} \right\}$$~~

3) Kulczynski =

$$\frac{1}{2} \left(\frac{S(A \cup B)}{S(A)} + \frac{S(A \cup B)}{S(B)} \right) = \frac{1}{2} (P(A|B) + P(B|A))$$

Used to find out imbalance ratio robust

4) Cosine (A, B) = $\frac{S(A \cup B)}{\sqrt{S(A) \times S(B)}} = \sqrt{P(A|B) \times P(B|A)}$

Null transaction \rightarrow It is a transaction that does not contain any of the itemsets being examined.

Null Variant \rightarrow If its value is free from influence of null transaction.

Correlation χ^2 and lift are not Null Variant

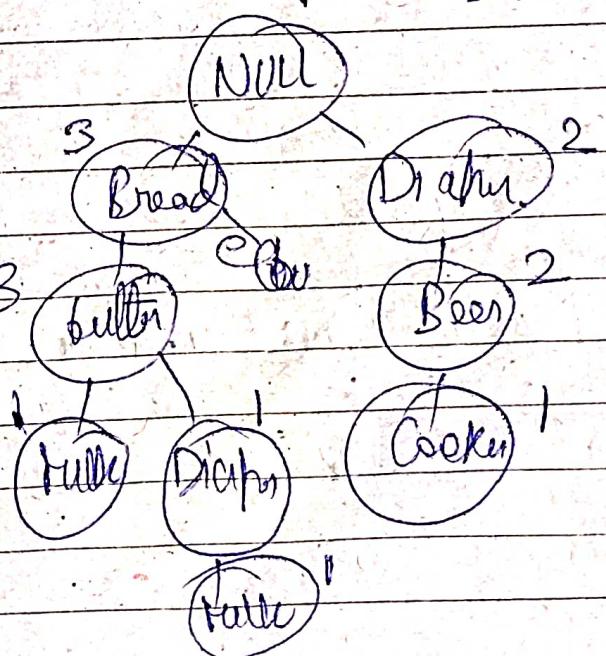
Pembalanck Ratio

$$DR(A, B) = \frac{|\text{Sup}(A) - \text{Sup}(B)|}{\text{Sup}(A) + \text{Sup}(B)}$$

$$\text{Sup}(A \cup B) - \text{Sup}(A) - \text{Sup}(B)$$

TID	Item Set	
1	{Bread, butter, Milk}	Bread, butter, Milk
2	{Bread, Butter}	Bread, Butter
3	{Beer, Cookies, Diapers}	Diapers, Beer, Cookies
4	{Milk, Diapers, Bread, Butter}	Bread, butter, Diapers, Milk
5	{Beer, Diapers}	Diapers, Beer

Item	Count	Records
Bread	3	1
butter	3	2
Milk	2	4
Beer	2	5
Cookies	1	6
Diapers	3	3



Chapt - 18

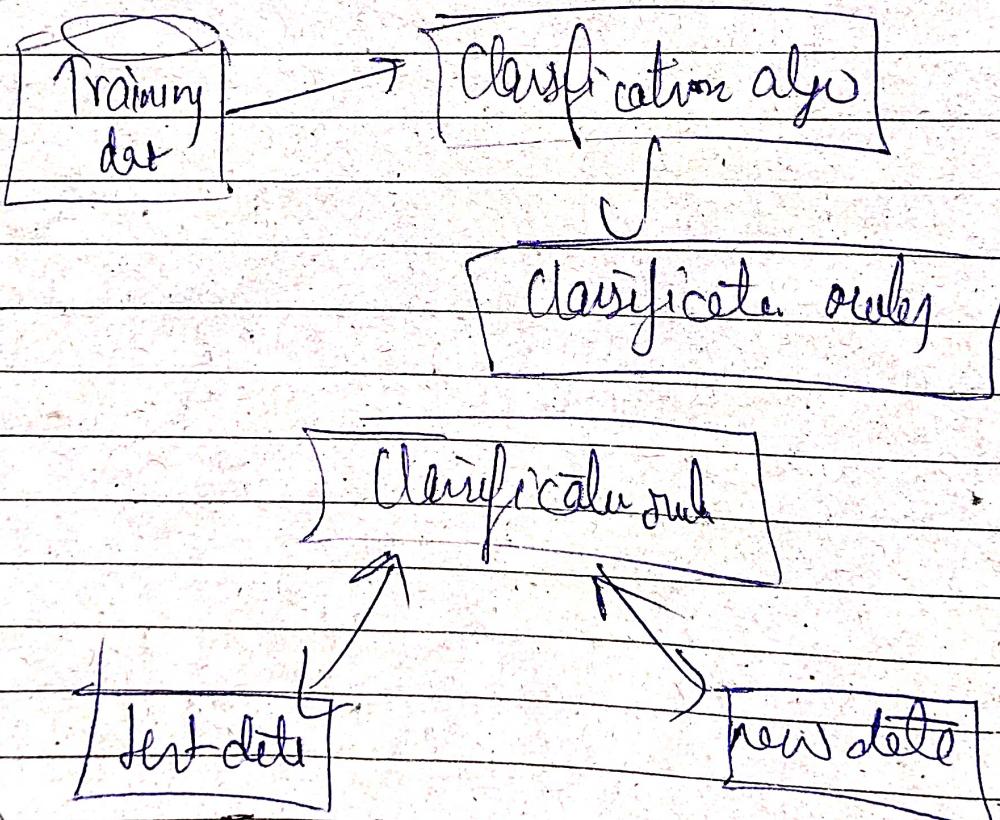
Classification

It is a form of data analysis that extracts model describing important data classes.

Data Classification \rightarrow (i) Learning step
 (ii) Classification step

a) Learning \rightarrow Training data are analyzed by a classification algo.

b) Classification \rightarrow Test data are used to estimate the accuracy of few classification rules



Decision Tree Induction:

It is the learning of decision tree from class-labeled training tuples.

Attribute Selection Measures

It is a heuristic for selecting the splitting criterion that "best" separate a given data partition, D , of class-labeled training tuples into individual classes.

Information Gain

ID3 uses infoGain as its attribute selection measure.

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2(p_i)$$

$\text{Info}(D)$ is known as Entropy

$$\text{Entropy} = - \sum_{i=1}^m \frac{p_i}{P+N} \log_2 \left(\frac{p_i}{P+N} \right)$$

$$= - \frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

$$\text{Info-Gain}(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \text{Entropy}(D)$$

$|D_j|$ acts as the weight of j th partition

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_{\text{Gain}}(D)$$

$$= \text{Entropy} - \text{Info}_{\text{Gain}}$$

tells how much would be gained by branching on A .

Entropy \rightarrow Uncertainty associated with given info

Info Gain \rightarrow Information Gained

ID3 algo.

Day	Weather	Temp	Humidity	Play
D1	Sunny	Hot	High	N
D2	Sunny	Hot	High	N
D3	Cloudy	Hot	High	Y
D4	Rainy	Mild	Normal	Y
D5	Rainy	Cool	Normal	Y

$$(1) \text{ Entropy(Entire)} \{3, -2\} = -\frac{P}{P+N} \log_{10}\left(\frac{P}{P+N}\right) - \frac{N}{P+N} \log_{10}\left(\frac{N}{P+N}\right)$$

$$= \frac{-3}{5} \log_{10}\left(\frac{3}{5}\right) - \frac{2}{5} \log_{10}\left(\frac{2}{5}\right)$$

$$= -\frac{P}{P+N} \log_{10}\left(\frac{P}{P+N}\right) - \frac{N}{P+N} \log_{10}\left(\frac{N}{P+N}\right)$$

$$\log_{10}(2)$$

$$= 3.321 \left(\frac{-3}{5} \log_{10}(0.6) - \frac{2}{5} \log_{10}(0.4) \right)$$

$$= 3.321 / (0.133 + 0.159)$$

$$= 0.969$$

$$(i) \text{ Entropy (Weather) } \{ \text{Sunny} \{ 0, -2 \} \} = 3.321 \left[-\frac{0}{2} \log \left(\frac{0}{2} \right) - \frac{2}{2} \log \left(\frac{2}{2} \right) \right] \\ = 0$$

$$\text{Entropy (Cloudy) } \{ 1, 0 \} = -\frac{1}{2} \log \left(\frac{1}{2} \right) - \frac{0}{2} \log \left(\frac{0}{2} \right) = 0$$

$$\text{Entropy (Rain)} \{ 2, 0 \} = 0$$

$$\text{Infraenum} = \text{Entropy (entire)} - \frac{2 \times 0}{5} - \frac{1 \times 0}{5} - \frac{2 \times 0}{5}$$

$$\text{Infraenum} = 0.969$$

(ii) Temp

$$\text{E}(Hot) \{ 1, -2 \} = 3.321 \left(-\frac{1}{3} \log \left(\frac{1}{3} \right) - \frac{2}{3} \log \left(\frac{2}{3} \right) \right) \\ = 3.321 (0.159 + 0.117) \\ = 0.9165$$

$$\text{E}(Mild) \{ 1, 0 \} = 0$$

$$\text{E}(Cold) \{ 1, 0 \} = 0$$

$$\text{Infraenum (Temp)} = 0.969 - \frac{3 \times 0.9165}{5} = 0.419$$

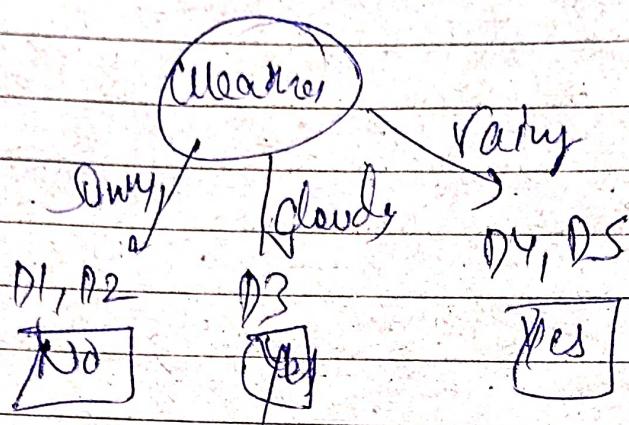
(iii) Humidity

$$\text{E}(High) \{ 1, -2 \} = 0.9165$$

$$\text{E}(Normal) \{ 2, 0 \} = 0$$

$$\text{Infraenum} = 0.419$$

Maximum Gini = 0.91



Gini Index

Purity of Data Set

$$Gini(D) = 1 - \sum_{i=1}^m (P_i)^2$$

$$Gini_A(D) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2)$$

$$\Delta Gini(A) = Gini(D) - Gini_A(D)$$

$$Gini_A(D) = \frac{10}{14} Gini(D_1) + \frac{4}{14} Gini(D_2)$$

Ex Day	Weather	Temb	Play
D1	Sunny	Hot	No
D2	Sunny	not	No
D3	Cloudy	not	Yes
D4	Rainy	Wet	Yes
D5	Rainy	Cool	No

$$\text{Gini(Entire)} \leftarrow \{3, -2\} = 1 - \sum p_i^2$$

$$= 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.48$$

Gini(Weather)

$$\text{Gini(Sunny)} \leftarrow \{0, -2\} = 1 - \left(\frac{0}{2}\right)^2 - \left(\frac{2}{2}\right)^2 = 0$$

$$\text{Gini(Cloudy)} \leftarrow \{1, 0\} = 0$$

$$\text{Gini(Rainy)} \leftarrow \{2, -0\} = 0$$

$$\text{Gini(Cilector)} = \frac{2 \times 0}{5} + \frac{1 \times 0}{5} + \frac{2 \times 0}{5} = 0$$

Gini(Temp)

$$\text{Gini(Hot)} \leftarrow \{1, -2\} = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 1 - 0.111 - 0.444 = 0.445$$

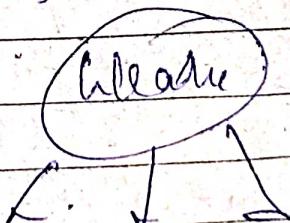
$$\text{Gini(Mild)} \leftarrow \{1, 0\} = 0$$

$$\text{Gini(Cool)} \leftarrow \{1, 0\} = 0$$

$$\text{Gini(Temp)} = \frac{3 \times 0.445 + 1 \times 0 + 0 \times 1}{5} = 0.267$$

Gini(Weather) < Gini(Temp)

by importance



Day	Weather	Temp	age	income	student	Grade	buy
D1	Sunny	Hot	Y	H	no	F	no
D2	Sunny	Hot	M	H	no	E	no
D3	Cloudy		M	H	no	F	Y
D4	Rainy		S	M	no	F	Y
D5	Rainy		S	L	yes	F	Y
D6	Rainy		S	L	yes	E	N
D7	Cloudy		M	L	yes	E	Y
D8	Sunny		Y	M	no	F	N
D9	Sunny		Y	L	yes	F	Y
D10	Rainy		S	M	yes	F	Y
D11	Sunny		Y	M	yes	E	Y
D12	Cloudy		M	M	no	E	Y
D13	Cloudy		M	H	yes	F	X
D14	Rainy		S	M	no	E	N

$$E(\text{Entire}) \{g, -S\} = \frac{-9}{14} \log\left(\frac{9}{14}\right) - \frac{5}{14} \log\left(\frac{5}{14}\right) = 0.34$$

~~Intercatgory Age~~

$$E(\text{Young}) \{g, -S\} = 0.92$$

$$E(\text{Middle}) \{g, -S\} = 0$$

$$E(\text{Senior}) \{g, -S\} = 0.92$$

$$IC_1 = 0.34 - \frac{5}{14} \times 0.92 - \frac{4}{14} \times 0 - \frac{5}{14} \times 0.97 = 0.24$$

Cart Algo

Age	Youth	5	Yes	2	
			No	3	→ greater
Middle	4		Yes	4	
			No	0	
Senior	5		Yes	3	
			No	2	

Attribute	Rules	Error (Mis/Total)	Total
Age	Youth → NO	2/5	
	Middle → Yes	0/4	4/4
	Senior → Yes	2/5	

Bank Income

High	4	Yes	2	Rules	Error	Total
		No	2	High → Yes/No	2/4	
Medium	6	Yes	4	Medium → Yes	2/6	5/14
		No	2	Low → Yes	1/4	
Low	4	Yes	3			
		No	1			

Blister

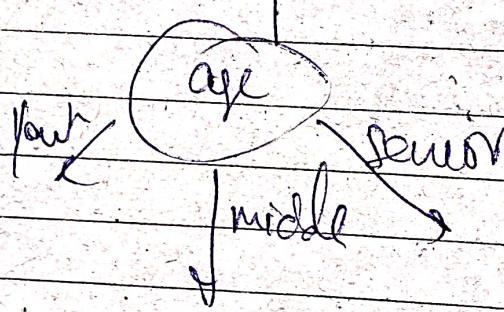
Yes	7	Yes	8	Rules	Error	
		No	1	Yes → Yes	1/7	8/14
NO	7	Yes	3	NO → NO	3/7	
		No	4			

Credit

			Rules	Error
Fair	8	Yes 6 No 2	Fair \rightarrow Yes Excellent \rightarrow Yes/No	2/8 9/10
Excellent	6	Yes 3 No 3		5/10

	Rules	Error	Total
Age \rightarrow	Young \rightarrow no Middle \rightarrow yes Senior \rightarrow yes	2/5 0/4 2/5	4/14
Student	Yes \rightarrow yes no \rightarrow no	1/1 2/2	4/4

Same
direct = 0



Tree Pruning

To Reduce Overfitting

Two approaches \rightarrow (i) Pre pruning
 (ii) Post pruning

Pre Pruning \rightarrow halting its construction early

Upon halting, the node becomes a leaf. The leaf may hold the most frequent class among the subset tuple or the probability distribution of those tuple.

When constructing a tree, measures such as statistical significance, information gain, Gini index, and so on, can be used to check the goodness of split.

If partitioning the data at a node would result in a split that falls below a prespecified threshold, then further partitioning is halted.

However choosing an appropriate threshold is very difficult.

High threshold \rightarrow oversimplified trees

Low threshold \rightarrow very little simplification

Various measure \rightarrow max-depth, max-leaf etc.

"Pre-Pruning is used on Big Data set"

Post Pruning \rightarrow remove subtrees from a fully grown tree.

A subtree at a given node is pruned by removing its branches and replacing it with a leaf. The leaf is labeled with the most frequent class among all subtrees being replaced.

"It is used on Small Data Set"

CART \rightarrow Cost Complexity pruning algo \rightarrow no of leaves error rate

C4.5 \rightarrow pessimistic pruning \rightarrow error rates

Bayesian Classifiers

Bayesian classifiers are statistical classifiers. They can predict class membership probability such as the probability that a given tuple belongs to a particular class.

→ Bayes' theorem

→ Naive Bayesian classifier

Bayes Theorem

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

$$P(A|B) \times P(B) = P(B|A)P(A)$$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

likelihood

H → hypothesis
 D → Date classifier

$$P(H|D)$$

posterior

$$P(D|H) \times P(H)$$

$$P(D)$$

prior

marginal

$P(\text{H}) \rightarrow \text{Prior} \rightarrow \text{Probability of hypothesis before giving evidence}$

Naive Bayes Classifier

Determine if Weather = Sunny, Temp = Cool, Humidity = High, Windy = Strong

$$P(\text{Yes}) = \frac{9}{14} = 0.64 \quad P(\text{No}) = \frac{5}{14} = 0.36$$

Conditional Probability of each feature

Weather	Yes	No	Humidity	Yes	No
Sunny	Sunny & Yes Total Yes	Sunny & NO Total NO	High	3/9	4/5
	= 2/9	3/5	Normal	6/9	1/5
Cloudy	4/9	0/5	Temp	Yes	No
Rainy	3/9	2/5	Hot	4/9	2/5
			Mild	4/9	2/5
			Cool	3/9	1/5

Windy	Yes	No
Strong	3/9	3/5
Weak	6/9	2/5

$$V_{NB} = \underset{V_j \in \{\text{Yes}, \text{No}\}}{\operatorname{argmax}} P(V_j) \cdot \prod_i P(a_i | V_j)$$

$$V_{NB}(\text{Yes}) = P(\text{Yes}) \cdot P(\text{Sunny} | \text{Yes}) \cdot P(\text{Cool} | \text{Yes}) \cdot P(\text{High} | \text{Yes}) \cdot P(\text{Strong} | \text{Yes})$$

$$= \frac{9}{14} \times \frac{2}{3} \times \frac{3}{5} \times \frac{3}{9} \times \frac{3}{9} = 0.053$$

$$V_{NB}(nw) = P(nw) \cdot P(\text{Sunny}/nw) \cdot P(\text{Cool}/nw) \cdot P(\text{High}/nw) \cdot P(\text{Snowy}/nw)$$

$$= \frac{3}{14} \times \frac{3}{5} \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5} = 0.026$$

To overcome Zero Probability Problem

Laplace Smoothing

$$P_k(x|y) = \frac{\text{Count}(x,y) + k}{\text{Count}(y) + k|X|}$$

$k \rightarrow$ Reference Smoothing parameter > 0

$|X| \rightarrow$ no of dimension / features \rightarrow weaker
 { Sun } \rightarrow
 { Rain } \rightarrow
 { Cloud } \rightarrow

$$\underset{x=1}{P}(\text{Outlook} = \text{Cloudy} | nw) = \frac{\text{Count}(\text{Outlook} = \text{Cloudy} | nw) + k}{\text{Count}(nw) + k|X|}$$

$$= \frac{1}{8}$$

Fruit = {Yellow, Sweet, Long}

Fruit	Yellow	Sweet	Long	Total
Orange	350	450	0	600
Banana	400	300	350	1050
Others	50	100	80	230
Total	800	850	400	1200

$$P(\text{Yellow} | \text{Orange}) = \frac{P(\text{Orange} | \text{Yellow}) \times P(\text{Yellow})}{P(\text{Orange})}$$

$$\frac{350}{800} \times \frac{800}{1200} = 0.55$$

$$P(\text{Sweet} | \text{Orange}) = P(\text{Orange} | \text{Sweet}) \times P(\text{Sweet}) = \frac{450}{850} \times \frac{850}{1200}$$

$$P(\text{Orange}) = \frac{600}{1200}$$

$$= 0.69$$

$$P(\text{Long} | \text{Orange}) = P(\text{Orange} | \text{Long}) \times \frac{P(\text{Long})}{P(\text{Orange})} = \frac{0}{1000} = 0$$

$$P(\text{Fruit} | \text{Orange}) = P(Y|O) P(S|O) \times P(L|O) = 0$$

Similarly

$$P(\text{Fruit} | \text{Banana}) = P(Y|B) P(S|B) P(L|B) = 1 \times 0.85 \times 0.9 = 0.765$$

$$P(\text{Fruit} | \text{Others}) = P(Y|\text{others}) P(S|\text{others}) P(L|\text{others}) = 0.12 \times 0.12 \times 0.12 = 0.001728$$

? greatest \rightarrow fruit = Banana

Rule Based Classification

Using IF - Then Rules

Rule antecedent
or
Precondition

Yule

Consequent

Rule Confliction

Conflict Resolution strategy to figure out which rule gets the fire and assign its class predictor to X.

(i) Size Ordering → assign the highest priority to the triggering rule that has the "toughest" requirements, where toughness is measured by the rule antecedent size.

(ii) The Rule Ordering →

a) Class based → classes are sorted in order of decreasing "importance" such as by decreasing order of prevalence.

b) Rule Based → Organized into a long priority list according to some measure of rule quality such as accuracy, coverage, etc.

$$\text{Coverage } (R) = \frac{n_{\text{covers}}}{|D|}$$

$$\text{Accuracy } (R) = \frac{n_{\text{correct}}}{n_{\text{covers}}}$$

Rule induction using a Sequential covering algo

Rule_Set = \emptyset

for each class c do

Rule = Learn_One_Rule(D, att_Value, c)

remove Rule covered by Rule from D

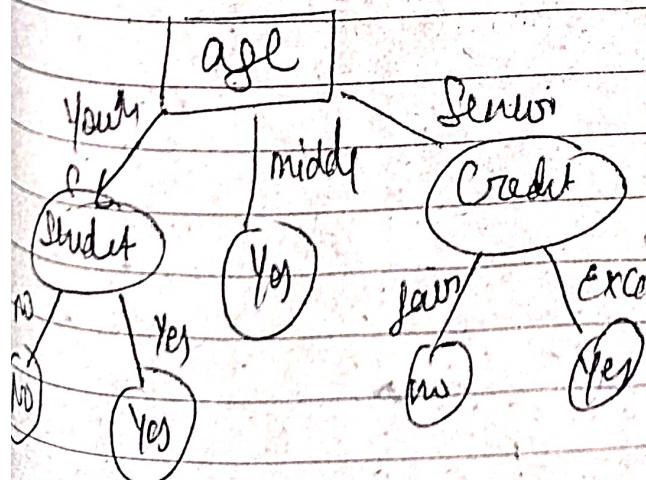
Rule_Set = Rule_Set + Rule;

Until terminality condition

end for

return Rule_Set;

Rule Extraction from Decision Tree



R1: If age = Young and student = no

then buy Computer = no

R2: If age = Young and student = Yes

then buy Computer = yes

R3: If age = middle and Fair

Computer = yes

R4: If age = senior and Credit = fair

then buy Computer = no

R5: If age = senior and Credit = excellent then buy Computer = yes

Characteristics of Rule

- ① Mutually Exclusive \rightarrow we cannot have rule conflict here because now the two rules will triggered for the same tuple.
- ② Exhaustive \rightarrow there is one rule per tuple, and any tuple has ~~one~~ \rightarrow for each possible attribute value combination.

Rule Duality Measure

FOIL \rightarrow first order induction learning

$$\text{① FOIL error} = \text{pos}' \left(\log_2 \left(\frac{\text{pos}'}{\text{pos}' + \text{neg}'} \right) - \log_2 \left(\frac{\text{pos}}{\text{pos} + \text{neg}} \right) \right)$$

$$\text{② Likelihood Ratio} = 2 \sum_{j=1}^m f_j \log \left(\frac{f_j}{e_j} \right)$$

Metrics for Evaluating Classifier Performance

		Predicted		Total
		Yes	No	
Actual	Yes	TP	FN	P
	No	FP	TN	N
Total	P	N	P+N	

Measure

Formula

accuracy, true positive rate

$$\frac{TP + TN}{P + N}$$

error rate, misclassified rate

$$\frac{FP + FN}{P + N}$$

sensitivity, true positive rate,

$$\frac{TP}{P}$$

Recall

$$P = \frac{TP}{(TP + FN)}$$

specificity, true negative rate,

$$\frac{TN}{N} = \frac{TN}{(FP + TN)}$$

Precision

$$\frac{TP}{TP + FP}$$

F, F₁, F₂ score

$$2 \times \text{Precision} \times \text{Recall}$$

harmonic mean of precision & recall

Precision + Recall

F_β; β → non-negative

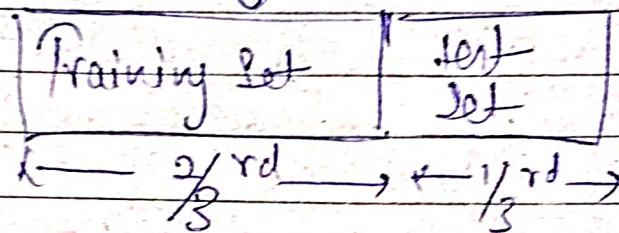
$$\frac{(1+\beta^2) \times \text{Precision} \times \text{Recall}}{\beta^2 \times \text{Precision} + \text{Recall}}$$

Holdout Method

The given data are randomly partitioned into two independent sets, a training set and test set.

[Data (Given)]

↓



Training set is used to derive the model.

Test set is used to estimate the model's accuracy.

The estimate is heuristic because only a portion of the initial data is used to derive the model.

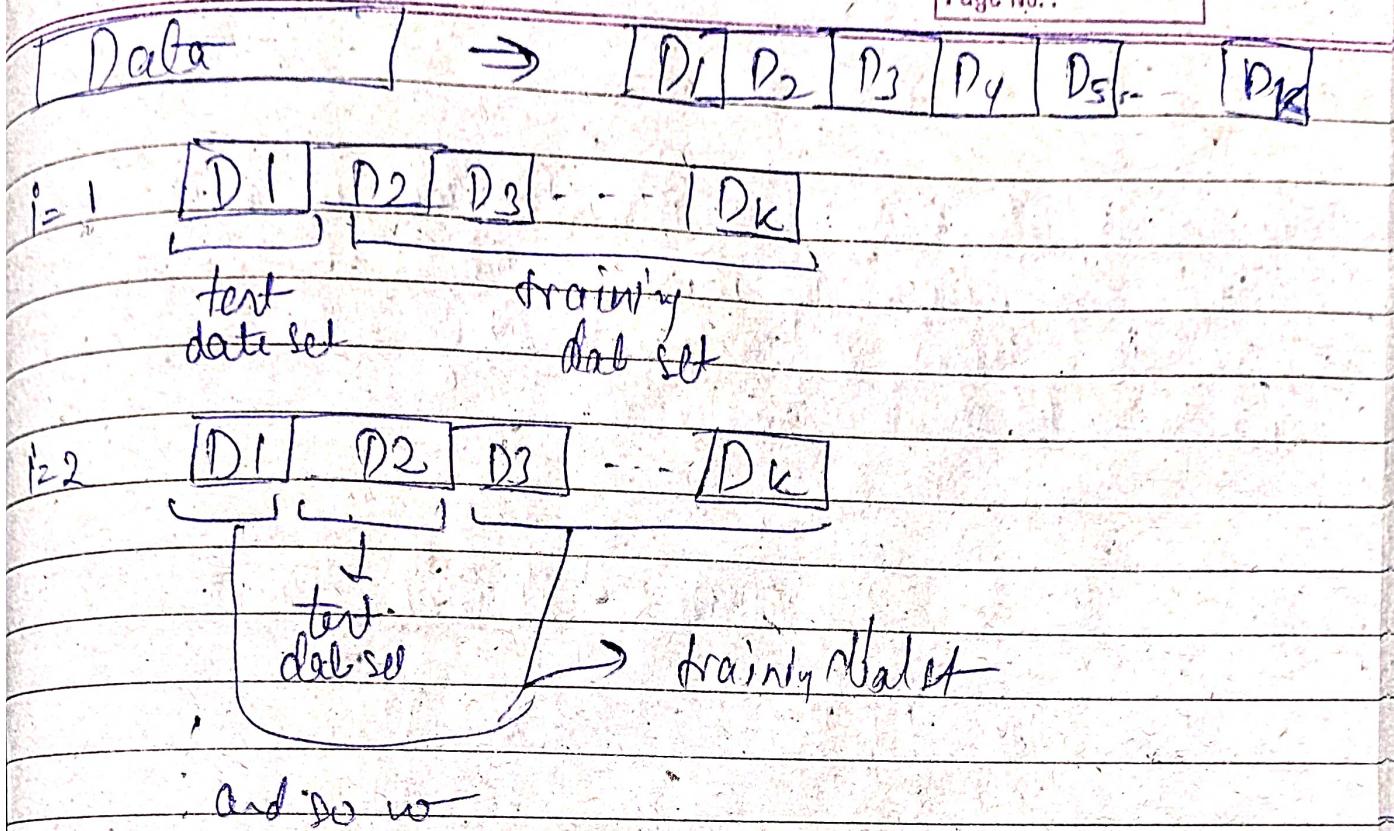
Random Subsampling

It is a variation of holdout method in which holdout method is repeated K times.

The overall accuracy \rightarrow average of each accuracy.

CROSS VALIDATION

K-FOLD Validation



Accuracy = Overall no of correct classifications / Total no of tuples in initial data

Leave-one-out

It is a special case of k -fold Cross Validation where k is set to the number of initial tuples.

Only one sample is "left out" at a time for the test set.

Stratified k -fold

The folds are stratified so that the class distribution of the tuples in each fold is approx the same as that in initial data.

Half-filled 10-fold Cross Validation \rightarrow accuracy becomes low bias and variance

Bootstrap

The Bootstrap method samples for given training tuples uniformly with replacement. That is each time a tuple is selected, it is equally likely to be selected again and re-added to the training set.

ROC Curve Receiver Operating Characteristic

TPR vs FPR

$$TPR = \frac{TP}{TP+FN}$$

\downarrow

TP

P

\downarrow

FP

N

$$FPR = \frac{FP}{FP+FN}$$

True Positive

$$P \rightarrow P$$

False Positive

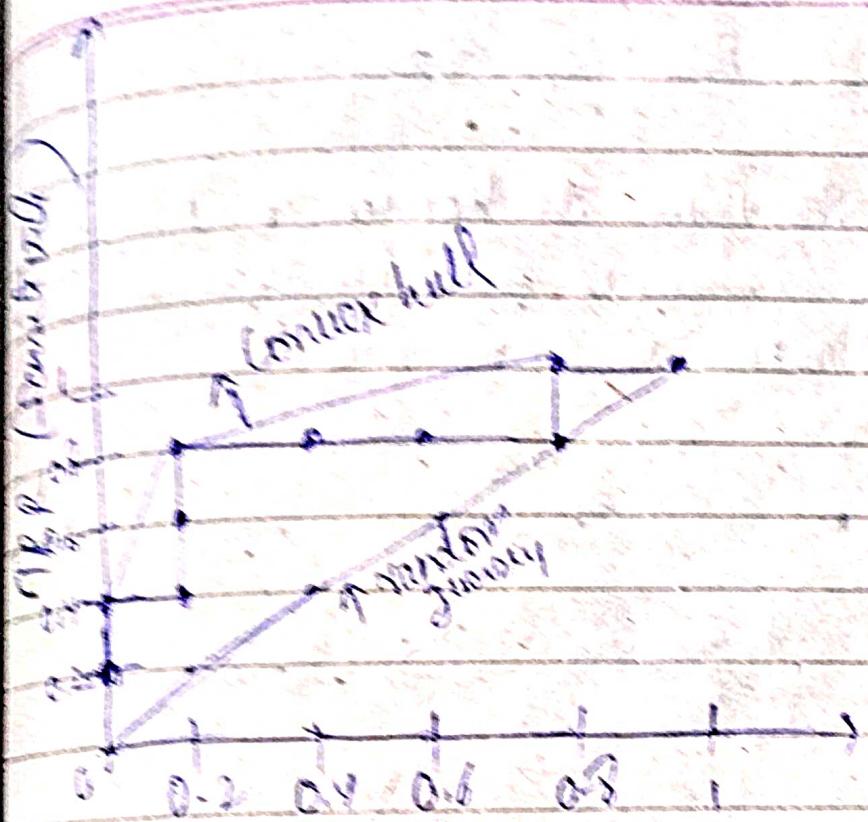
$$N \rightarrow P$$

Total Positive

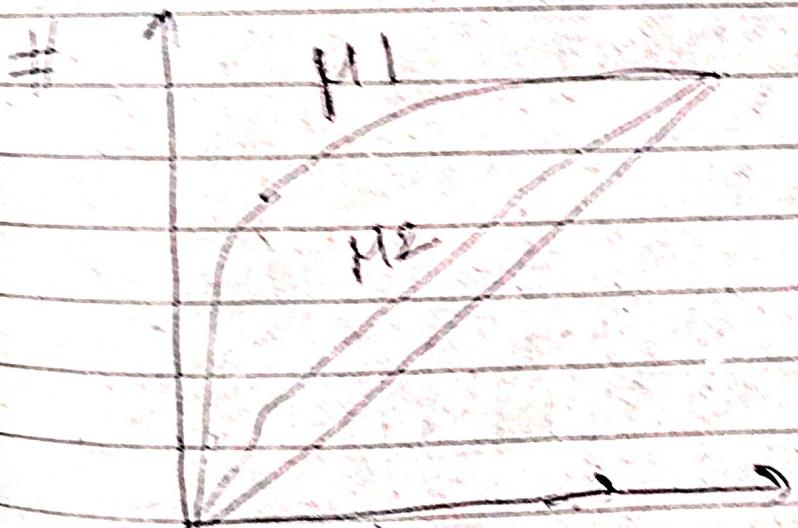
Total Negative

$$\frac{\text{Total } P}{\text{Total } N}$$

Ex. Tuples #	Class	Prob	TP	FP	$TPR = \frac{TP}{S}$	$FPR = \frac{FP}{S}$
1	$P \rightarrow P$	0.90	1	0	$\frac{1}{5} = 0.2$	0
2	$P \rightarrow P$	0.80	1+1=2	0	0.4	0
3	$P \rightarrow N$	0.80	2	0+1=1	0.4	0.2
4	$P \rightarrow P$	0.60	2+1=3	1	0.6	0.2
5	$P \rightarrow P$	0.55	3+1=4	1	0.8	0.2
6	$P \rightarrow N$	0.54	4	1+1=2	0.7	0.4
7	$P \rightarrow N$	0.53	4	2+1=3	0.8	0.6
8	$P \rightarrow N$	0.51	4	3+1=4	0.8	0.3
9	$P \rightarrow P$	0.50	2+1=3	4	1	0.8
10	$P \rightarrow N$	0.40	5	4+1=5	1	1



FPR ($1 - \text{Specificity}$)



The ROC Curve is to the diagonal line, the less accurate the model is.

Thus H1 is more accurate