

3) Partial materialization:

Selectively compute a proper subset of the whole set of possible cuboids if consider three factors

i) Identify the subset of cuboids or subcubes to materialize.

ii) Use the materialized cuboids during query processing cuboid during load and refresh.

Association Rule:

Frequent + pattern.

→ support → 2% mobile → Tempered glass

→ confidence → 60%

Range:

i) min-max nor : [new-min, new-max]

ii) z-score nor : [-∞, ∞]

iii) " using mean absolute deviation : [0, ∞]

iv) Decimal scaling : [-1 to 1]

v) Normalization : [0, 1] & [-1 + 1]



Association Rules:

$A \rightarrow B$

I_1	I_2	I_3	I_n
1	0	0	1	0 0 0 0

Support = 2% confidence = 60%

Let, $I = \{I_1, I_2, \dots, I_m\}$ be a set of items.

$T \rightarrow$ Transaction

⇒ let 'A' be a set of items in 'T' is said to contain A if and only if $A \subseteq T$

⇒ An association rule $A \rightarrow B$, where $ACI, BCI \neq \emptyset$ let $A \rightarrow B$ holds in transaction set (T)

support ($A \rightarrow B$) = $P(A \cup B)$

confidence ($A \rightarrow B$) = $P(B|A)$

confidence = $\frac{\text{support}(A \cup B)}{\text{support}(A)}$

Relative support = $\frac{\text{support}(A \cup B)}{\text{support count}(A \cup B)}$

closed frequent itemset = $\frac{\text{support count}(A)}$

* An item set X is closed in dataset S if there exists no proper super itemset Y such that Y has the same support count as X in S .

* An item set X is maximal if there exists no super itemset Y such that $X \subset Y$ and Y is frequent in S . (IA on site)

Q) Suppose that the database has only 2 transactions.

{ $(a_1, a_2, a_3, \dots, a_{100})$, $(a_1, a_2, a_3, a_4, \dots, a_{50})$ }

check solution

Multi-level Association Rules:

→ suppose ' $x!$ ' is a variable, representing a customer buys (' x ', "Computer") → buys($x!$, "HP pointer")

buys (' x ', "Laptop computer") → buys (' $x!$ ', "HP pointer")

In this multi-level association rules, rules are referred at different levels of association.

→ If the items are attribute in an association rule refer only one dimension that it is single dimension

If the rule refers two or more dimension such as age, income, buys etc then it is multi-D association rule.

Apriori Algorithm:

→ 1994

to implement frequent item sets in a good manner

R. Agrawal & R. Srikant

- It uses level-wise search, where K items are used to explore $K+1$ item set.
 - It is used for mining boolean frequent itemsets.
 - It uses prior knowledge of frequent items.
 - First, the set of frequent 1-itemset is found by scanning the database to accumulate count for each item and collecting those items that satisfy minimum support. The resulting set is denoted by L_1 .
 - Next, L_1 is used to find L_2 and so on.
 - Finding each L_k requires a full scan of database.
- Apriori property:
- All non-empty subsets of a frequent itemset must also be frequent.
 - If an item A is added to the itemset I , then the resulting itemset $I \cup A$ cannot occur more frequently than I .
 - This property belongs to special category of property called Antimonotone in the sense that if the set cannot pass a test all of its super set will fail in the same test as well.
 - ⇒ Apriori L_{k-1} is used to find L_k for $k \geq 2$. It is a two-step process followed by joined or prune actions.
 - 1) Join step: To find L_k a set of candidates k -itemsets is generated by joining L_{k-1} with itself. The candidates are denoted by C_k .
 - 2) Prune step: C_k is a super set of L_k . A scan of database to determine sub-count of each candidate in C_k would result in the determination of L_k .

C_k can be used, to reduce the size of C_k apriori property is used. Only C_{k-1} itemset where not frequent it can't be a subset of frequent itemset. (ex. $\{I_1, I_2\}$ is not frequent)

- Ex: $T_1 \rightarrow I_1, I_2, I_5$ frequent itemset. I_5 is the residue itemset
 $T_2 \rightarrow I_2, I_4$ residue itemset
 $T_3 \rightarrow I_2, I_3$ residue itemset - more prove left
 $T_4 \rightarrow I_1, I_2, I_4$ residue itemset $(C_2 - 2) \leftarrow 2$ items
 $T_5 \rightarrow I_1, I_3$ frequent itemset
 $T_6 \rightarrow I_2, I_3$ residue itemset
 $T_7 \rightarrow I_1, I_3$ frequent itemset
 $T_8 \rightarrow I_1, I_2, I_3, I_5$ frequent itemset with support 4
 $T_9 \rightarrow I_1, I_2, I_3$ frequent itemset $\{I_1, I_2, I_3\} = C_3$

C_1 is first frequent itemset

Support for count of each candidate

I_1	6
I_2	7
I_3	6
I_4	2
I_5	2

Support-count with minimum support count

Support = $\frac{\text{count}}{\text{total}} \times 100$

Compare I_1, I_2, I_3 with I_1, I_3, I_4

I_1, I_2	4
I_1, I_3	4
I_1, I_4	1
I_1, I_5	2
I_2, I_3	4
I_2, I_4	2
I_2, I_5	2
I_3, I_4	0
I_3, I_5	1
I_4, I_5	0

Generate C_2 from L_1

form L_2 from C_2

Support for count of each candidate

I_1, I_2, I_3	2
I_1, I_2, I_5	2
I_1, I_2, I_4	1
I_1, I_3, I_5	1
I_1, I_3, I_4	0
I_1, I_5, I_3	1
I_1, I_5, I_4	0
I_2, I_3, I_4	0
I_2, I_3, I_5	1
I_2, I_4, I_5	0

Support-count with minimum support count

Support = $\frac{\text{count}}{\text{total}} \times 100$

Compare I_1, I_2, I_3 with I_1, I_3, I_5

I_1, I_2	4
I_1, I_3	4
I_1, I_5	2
I_2, I_3	4
I_2, I_4	2
I_2, I_5	2

Generate C_3 from L_2

form L_3 from C_3

Support = $\frac{\text{count}}{\text{total}} \times 100$

I_1, I_2, I_3	2
I_1, I_2, I_5	2
I_1, I_3, I_5	2

* Generating rules from frequent itemset

$$\text{Confidence } (A \rightarrow B) = \frac{\text{support count } (A \cup B)}{\text{support - count } (A)}$$

- * For each frequent itemset 'l', generate all non-empty subsets of l
* For every non-empty subset 's' of 'l' output the rule $s \rightarrow (l-s)$ where $\frac{\text{support - count } (l)}{\text{support - count } (s)}$ min. confidence.

Suppose the data contains frequent itemset

$l = \{I_1, I_2, I_5\}$ then subset of l are $\{I_1, I_2\}$,
 $\{I_1, I_5\}$, $\{I_2, I_5\}$, $\{I_1\}$, $\{I_2\}$, $\{I_5\}$.

$I_1 \wedge I_2 \rightarrow I_5$ conf. = $\frac{2}{4} = 50\%$.

$I_1 \wedge I_5 \rightarrow I_2$ conf. = $\frac{2}{2} = 100\%$.

(FP Growth (FP-Tree))

Apriori Algo. suffer from main problems:

1) It may need to generate a huge no. of candidate set

2) It may need to repeatedly scan the database

FP Growth : (FP-Tree)

In this the first scan of the DB is same as apriori which derives the set of frequent items. and they support_count

* Let the min support_count is "2"

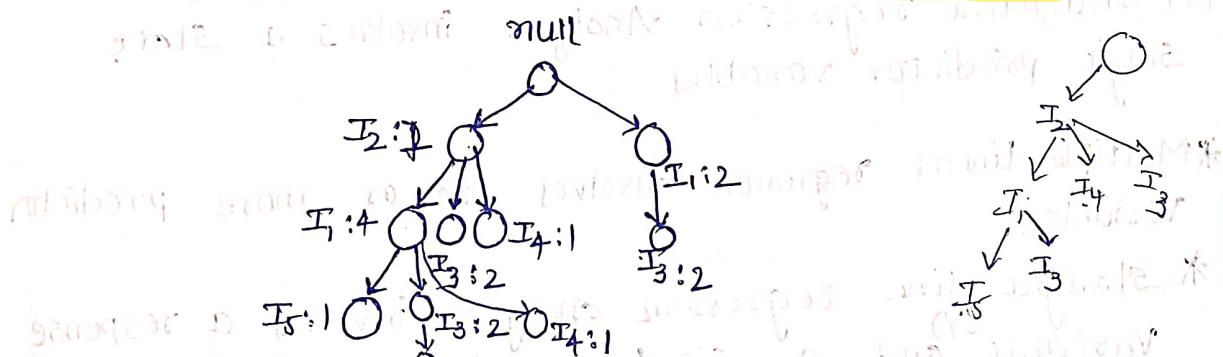
* The set of frequent items are sorted in the descending order of support_count, the resulting set or list is denoted by 'L' and thus

<u>Tid</u>	<u>Itemsets</u>
T_1	$I_1, I_2, I_5 \rightarrow I_2, I_1, I_5$
T_2	$I_2, I_4 \rightarrow I_2, I_4$
T_3	$I_2, I_3 \rightarrow I_2, I_3$
T_4	$I_1, I_2, I_4 \rightarrow I_2, I_1, I_4$
T_5	$I_1, I_3 \rightarrow I_1, I_3$
T_6	$I_2, I_3 \rightarrow I_2, I_3$
T_7	$I_1, I_3 \rightarrow I_1, I_3$
T_8	$I_1, I_2, I_3, I_5 \rightarrow I_2, I_1, I_3, I_5$
T_9	$I_1, I_2, I_3 \rightarrow I_2, I_1, I_3$

$$L = \{ \{I_2:7\}, \{I_1:6\}, \{I_3:6\}, \{I_4:2\}, \{I_5:2\} \}$$

* An FP tree is constructed as follows.

- 1) Firstly, create root of the tree labelled with null.
- 2) Scan the DB 'D', for the second time the items in each transaction are processed, in the 'L' order i.e., sorted acc to decreasing support term and the branch is created for each transaction.



- 3) An FP tree is mined as follows

- ⇒ Start from each frequent one. pattern (as an initial suffix pattern), construct its conditional pattern base then construct its conditional FP-tree and perform mining recursively in such a tree.
- ⇒ The pattern growth is achieved by concatenation of the suffix pattern with the frequent pattern generated from the conditional FP-tree.

Item	Conditional pattern Base	Conditional FP tree	frequent pattern generate
I ₅	{I ₂ , I ₁ : 1} {I ₂ , I ₁ , I ₃ : 1}	{I ₂ : 2, I ₁ : 2}	{I ₂ , I ₅ : 2}
I ₄	{I ₂ : 1} {I ₂ , I ₁ : 1}	{I ₂ : 2}	{I ₂ , I ₄ : 2}
I ₃	{I ₂ , I ₁ : 2} {I ₂ : 2} {I ₁ : 2} {I ₂ : 4, I ₁ : 2}		{I ₂ , I ₃ : 4}
I ₁	{I ₂ : 4}	{I ₂ : 4}	{I ₂ , I ₃ : 4}

Regression:

- * Regression can be used to model the relationship b/w one or more independent or predictor variables and a dependent or response variables predictor variables describes the tuples (making up the attribute vector).
- * In general the values of predictor variables are known Response variable is what we want to predict
- * A straightline regression Analysis involves a single predictor variables
- * Multiple linear regression involves one or more predictor variables.
- * Straight line regression analysis involves a response variable and a single predictor variable (x) it models y as a linear function of x .

$$y = b + w_1x$$
 where b, w_1 are regression co-efficients.
- * Regression coefficients can be thought of as a weight to go coe can equivalently write

$$y = w_0 + w_1x$$
- * These co-efficients can be solved by the method of least square, which estimates the best fit line which minimises the error between the actual data and the estimated of line.

* Regression co-efficient can be estimated by

$$\Rightarrow w_1 = \frac{\sum_{i=1}^{10} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{10} (x_i - \bar{x})^2}$$

$$w_0 = \bar{y} - w_1 \bar{x}$$

Write on desk

Year Experience	Salary (Rs)
3	30
8	57
9	64
13	72
10	36
3	43
6	59
11	90
21	20
1	83
16	

Find the Salary of person with 10 year of experience.

$$\bar{x} = \frac{3 + 8 + 9 + 13 + 3 + 6 + 11 + 21 + 1 + 16}{10} = 9.1$$

$$\bar{y} = \frac{30 + 57 + 64 + 72 + 36 + 43 + 59 + 90 + 20 + 83}{10} = 55.4$$

$$w_1 = \frac{(3-9.1)(30-47.9) + (8-9.1)(57-47.9) + (9-9.1)(64-47.9) + (13-9.1)(72-47.9) + (3-9.1)(36-47.9) + (1-9.1)(43-47.9) + (6-9.1)(59-47.9) + (11-9.1)(54-47.9) + (21-9.1)(90-47.9) + (1-9.1)(20-47.9) + (16-9.1)(83-47.9)}{(3-9.1)^2 + (8-9.1)^2 + (9-9.1)^2 + (13-9.1)^2 + (3-9.1)^2 + (1-9.1)^2 + (6-9.1)^2 + (11-9.1)^2 + (21-9.1)^2 + (1-9.1)^2 + (16-9.1)^2}$$

$$w_1 = 3.5$$

$$\Rightarrow w_0 = \bar{y} - w_1 x$$

$$= 55.4 - (3.5)(9.1)$$

$$= 55.4 - (3.5 \times 9.1)$$

$$= 23.5$$

$$y = 23.5 + 55.4(10)$$

$$= 23.5 + 554$$

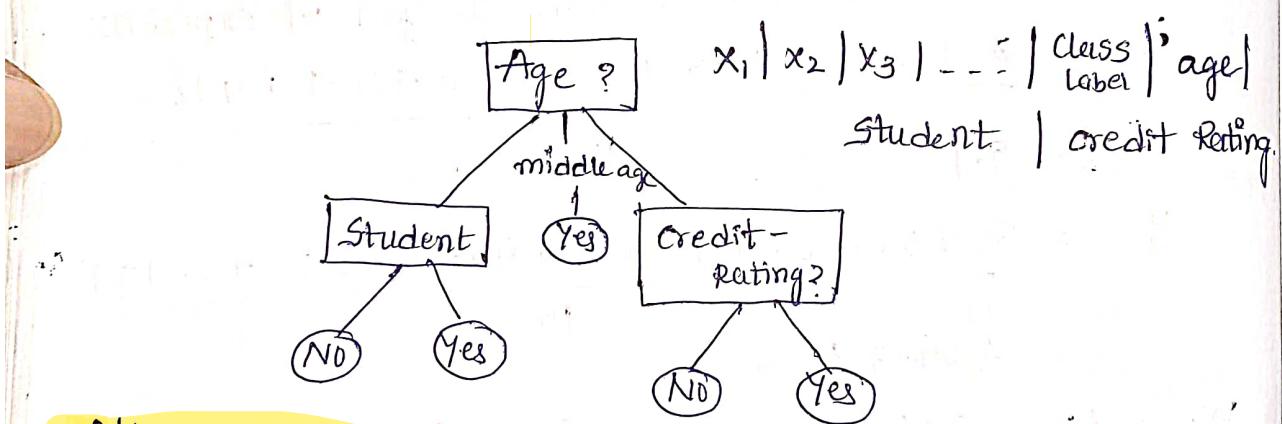
$$= 577.5$$

$$\begin{array}{r} 554 \\ 23.5 \\ \hline 577.5 \end{array}$$

Decision tree:

A decision tree is a flowchart by a tree structure where each internal node denotes a test or attribute, each branch represents an outcome of test and each leaf node (terminal node) holds a class label. The top most node in the tree is a root node.

Internal nodes represented by rectangles
 Leaf nodes " " ovals.



Attribute selection methods:

- ↳ Information gain
- ↳ Gain ratio
- ↳ Gini Index

- * Attribute selection measure determines how the tuples at a given node are to be split
- * The attribute having the best score for measure is chosen as the splitting attribute for the given tuple.

Information Gain:

- * The node with the highest info gain is chosen as the splitting attribute for node (N)
- * Expected information needed to classify a tuple is given by

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

where p_i = probability that an arbitrary tuple belongs to class C_i

p_i is estimated by

$ C_i, D $
$ D $

Input D is also known Info.

as Entropy of D

Write on desk

$$* \text{Info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{Info}(D_j)$$

where $|D_j|$ is estimated as weight of Jth position

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

(Information gain)

Age (U), Income (high) student (NO) credit rating class (NO) buys computer

Age	Income	Student	Credit rating	Class	Buy-computer
Youth	High	NO	Fair	NO	NO
Youth	High	NO	Excellent	NO	NO
Middle	High	NO	Excellent	NO	NO
Senior	Medium	NO	Fair	Yes	Yes
Senior	Low	Yes	Fair	Yes	Yes
Senior	Low	Yes	Fair	Yes	Yes
Middle	Low	Yes	Excellent	NO	NO
Youth	Medium	NO	Fair	Yes	Yes
Youth	Low	Yes	Fair	NO	NO
Senior	Medium	Yes	Fair	Yes	Yes
Youth	Medium	Yes	Fair	Yes	Yes
Middle	Medium	Yes	Excellent	Yes	Yes
Middle	High	NO	Excellent	Yes	Yes
Middle	High	Yes	Fair	Yes	Yes

Senior	medium	NO	Excellent	YES NO
--------	--------	----	-----------	--------

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

$$Info_{age}(D) = \frac{5}{14} \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}\right) + \frac{4}{14} \left(-\frac{4}{4} \log_2 \frac{0}{4}\right) + \frac{5}{14} \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}\right)$$

$$\approx 0.694$$

$$\begin{aligned} Gain(Age) &= Info(D) - Info_{age}(D) \\ &= (-0.694 + 0.94) \\ &= 0.246 \end{aligned}$$

$$Gain(Income) = 0.029$$

Gain ratio:

It is an extension of Information gain it attempts to overcome this bias it applies a kind of normalization to info gain using split information.

$$\Rightarrow \text{split info}_A(D) = \sum_{j=1}^v \frac{D_j}{D} \log_2 \frac{|D_j|}{|D|}$$

Gain ratio is defined as

minus

write on desk

$$\Rightarrow \text{Gain ratio}(A) = \frac{\text{Gain}(A)}{\text{split info}(A)}$$

The attribute with the max. gain ratio is selected as the splitting attribute

$$\begin{aligned} \text{split info}(D) &= -\frac{4}{14} \log_2 \frac{4}{14} - \frac{6}{14} \log_2 \frac{6}{14} - \frac{4}{14} \log_2 \frac{4}{14} \\ \text{Income} &\approx 1.57 \end{aligned}$$

$$\text{Gain Income} = 0.029.$$

$$\text{Gain ratio (Income)} = \frac{0.029}{1.58}$$

Gini Index:

* It is used in CART, gini index measures the impurity of data

$$Gini(D) = 1 - \sum_{i=1}^v p_i^2$$

where, P_i is the probability that a tuple in D belongs to the class C_i ; it is estimated by $\frac{c_i}{D}$

* The gini index considers a binary split of each attribute if a binary split on attribute A partitions D into D_1 and D_2 . The gini index of D is calculated as

$$\Rightarrow \text{Gini}_A(D) = \frac{|D_1|}{D} \text{Gini}_A(D_1) + \frac{|D_2|}{D} \text{Gini}_A(D_2)$$

The reduction in impurity is calculated as

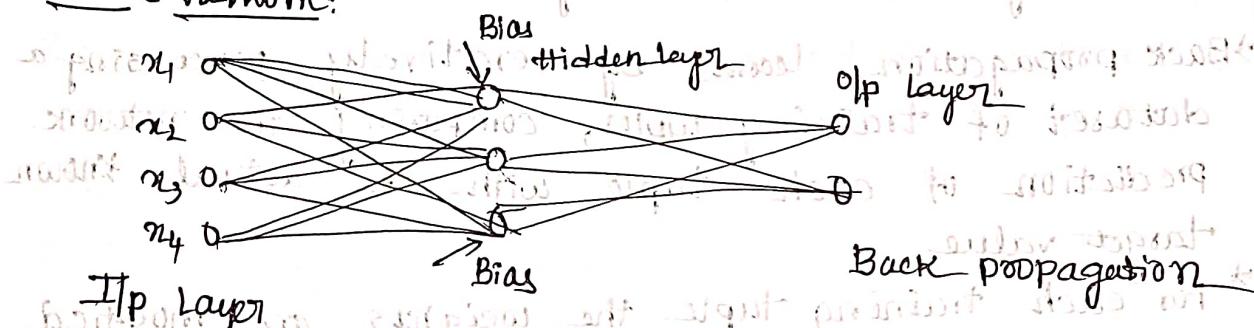
$$\Rightarrow \Delta \text{Gini}(A) = \text{Gini}(D) - \text{Gini}_A(D)$$

$$\text{Gini}(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2$$

$$\text{Gini}_A(D) = \frac{9}{14} \left(1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2\right) + \frac{6}{14} \left(1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2\right)$$

$$\approx 0.428$$

Neural network:



\Rightarrow NN is a set of connected input output units, in which each connection has a weight associated with it

\Rightarrow During the learning phase the network learns by adjusting weight associated with it so as to predict the correct class labels of input labels.

\Rightarrow Neural network layer is also referred to as Connectionist learning due to the connection units. Neural network involves forward pass and backward pass.

⇒ NN involves long training time and are therefore more suitable for applications where this is feasible.

- ⇒ The most popular Neural network Algo is back propagation.
- ⇒ The back propagation Algo. performs learning on a multi layer feed forward neural network.
- ⇒ It iteratively learns a set of weights for prediction of class label of tuples.
- ⇒ A multi-layer feed forward NN consists of an input layer one or more hidden layer and an output layer.
- ⇒ Each layer is made up of units the input to the network corresponds to attribute measured for each training tuple.
- ⇒ The Input are fed simultaneously into the units making up the input layer.
- ⇒ Back propagation, learns by iteratively processing a dataset of training tuples, comparing the network prediction of each tuple with the actual known target value.
- ⇒ For each training tuple the weights are modified so as to minimize the mean-squared error b/w the network prediction and the actual known target value.
- ⇒ These modifications are made in the backward direction i.e from the O/p layer through each hidden layer upto the 1st hidden layer, hence it is called back propagation.

$$I_j = \sum w_{ij} \cdot o_i + \theta_j$$

w_{ij} = weight of connection from unit i previous layer to unit j .

o_i = Output of the unit i from the previous layer

θ_j = Bias of the unit

used to vary the activity of unit

→ Each unit in the hidden layer and the o/p layer takes its net input and then apply the activation function to it.

→ The sigmoid func^c is generally used for this purpose given the ^{net} input i_j to unit j then O_j the output of unit j is computed as

$$O_j = \frac{1}{1 + e^{-i_j}}$$

→ This function also referred to as squashing function because it maps large into a smaller range (from 0 to 1).

Error is calculated as follows

Last layer $\text{Err}_j = O_j(1 - O_j)(T_j - O_j)$

T_j = known target value.

Hidden layer $\text{Err}_j = O_j(1 - O_j) \sum w_{jk} \text{Err}_k$

w_{jk} = weight of the connection from unit j to k

→ The weights and bias are updated to reflect the propagated error. weights are updated by the following equation.

$$\Delta w_{ij} = l \text{Err}_j O_i$$

$$w_{ij}^{\text{new}} = w_{ij} - \Delta w_{ij}$$

where l = learning rate that technical varies (0.01 to 1.0)

→ Bias are updated by the following equation

$$\Delta \theta_j = (l) \text{Err}_j$$

$$\theta_{j,\text{new}} = \theta_j - \Delta \theta_j$$

Chain Rule:

$$y = f(x) = \frac{1}{1 + e^{-x}}$$

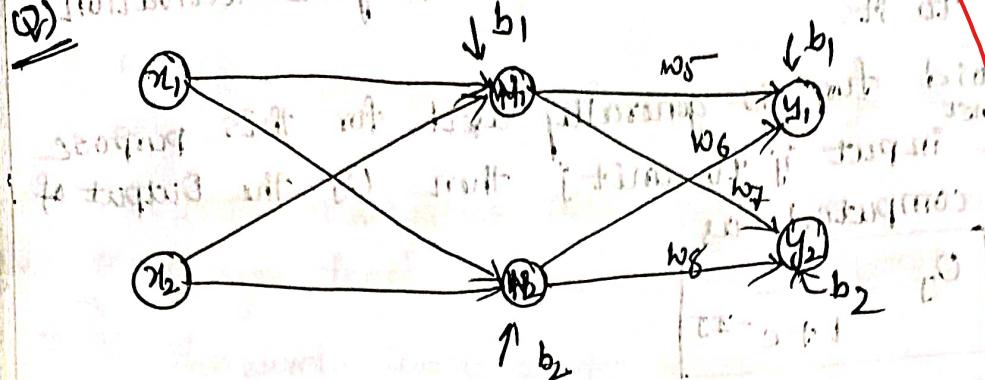
$$\frac{dy}{dx} = \left[\frac{1 - e^{-x}}{(1 + e^{-x})^2} \right] = \frac{e^{-x}}{(1 + e^{-x})^2}$$

$$f'(x) = \frac{e^{-x}}{(1 + e^{-x})^2}$$

$$f(x) \cdot (1 - f(x)) = \frac{e^{-x}}{(1 + e^{-x})^2}$$

$$f'(x) = f(x)(1 - f(x))$$

(Q)



$$x_1 = 0.05, b_1 = 0.35, T_1 \approx 0.0$$

$$x_2 = 0.10, b_2 = 0.6, T_2 = 0.99$$

$$w_1 = 0.45$$

$$w_5 = 0.4$$

$$w_2 = 0.2$$

$$w_6 = 0.45$$

$$w_3 = 0.25$$

$$w_7 = 0.5$$

$$w_4 = 0.3$$

$$w_8 = 0.55$$

$$H_1(\text{in}) = x_1 w_1 + x_2 w_2 + b_1 \approx 0.377$$

$$H_1(\text{out}) = \frac{1}{1 + e^{-H_1(\text{in})}} \approx 0.5932$$

$$H_2 = 0.596$$

$$y_1(\text{in}) = w_5 H_1(\text{out}) + w_8 H_2(\text{out}) + b_2$$

$$y_1(\text{out}) = \frac{1}{1 + e^{-y_1(\text{in})}}$$

$$E_{\text{total}} = \frac{1}{2} \sum (\text{Target} - \text{Actual o/p})^2$$

$$= \frac{1}{2} (T_1 - y_1(\text{out}))^2 + \frac{1}{2} (T_2 - y_2(\text{out}))^2$$

$$E_{\text{total}} = E_1 + E_2$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = \frac{\partial E_1}{\partial w_5} + \frac{\partial E_2}{\partial w_5}$$

$$\frac{\partial E_1}{\partial w_5} = \frac{\partial E_1}{\partial y_1(\text{out})} \times \frac{\partial y_1(\text{out})}{\partial y_1(\text{in})} \times \frac{\partial y_1(\text{in})}{\partial w_5}$$

$$= (y_1(\text{out}) - T_1) [y_1(\text{out})(1 - y_1(\text{out}))]$$

$$\frac{\partial E_{\text{total}}}{\partial w_1} = \frac{\partial E_1}{\partial w_1} + \frac{\partial E_2}{\partial w_2}$$

$$\frac{\partial E_1}{\partial w_1} = \frac{\partial E_1}{\partial y_1(\text{out})} \times$$

$$\frac{\partial y_1(\text{out})}{\partial y_1(\text{in})} \times \frac{\partial y_1(\text{in})}{\partial H_1(\text{out})} \times \frac{\partial H_1(\text{out})}{\partial H_1(\text{in})} \times \frac{\partial H_1(\text{in})}{\partial w_1}$$

Assign

[Diff b/w OLAP and OLTP?]

Bayes classifier:

It is a Supervised classification technique. The model assumes all input attributes are of equal importance and independent of one another.

Bayes classifier is based on the Bayes theorem, which is stated as

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

H is a hypothesis to be tested

E is the evidence associated with H

$P(\frac{A}{B})$ it is the probability of A given that B already occurs $P(B) > 0$.

$P(H)$ = is a prior probability which denotes the prob of hypothesis before presentation of any evidence.

Magazine promotion	watch promotion	Life Insurance	credit card	Gender class
Yes	NO	NO	NO	M
Yes	YES	YES	YES	F
NO	NO	NO	NO	M
Yes	YES	YES	YES	M
Yes	NO	NO	NO	F
NO	NO	NO	NO	F
Yes	YES	NO	NO	F
NO	NO	NO	YES	M
Yes	NO	NO	NO	M
Yes	NO	NO	NO	M
Yes	YES	NO	NO	M
		YES	NO	F

which the above given data find gender of tuple where magazine $\rightarrow S$, watch $\rightarrow N$, life insurance $\rightarrow N$, credit card $\rightarrow N$, gender $\rightarrow M$ or F.

$$\frac{P(\text{Gender} = \text{Male})}{E} = \frac{P(E | \text{Gender} = \text{Male}) P(\text{Gender} = \text{male})}{P(E)}$$

$$\Rightarrow P(MP = \text{Yes} | \text{Gender} = \text{male}) = \frac{4}{6}$$

$$\Rightarrow P(WP = \text{Yes} | \text{Gender} = \text{male}) = \frac{2}{6}$$

$$\Rightarrow P(CL = \text{No} | G = M) = \frac{4}{6}$$

$$\Rightarrow P(CC = \text{No} | G = M) = \frac{4}{6}$$

$$\left(\frac{P(G = M)}{E} \right) = \frac{4}{6} \times \frac{2}{6} \times \frac{4}{6} \times \frac{4}{6} \times \frac{6}{10} = 0.059$$

$$\Rightarrow \left(\frac{P(\text{Gender} = \text{female})}{E} \right) = \frac{P(E | G = F) P(G = F)}{P(E)}$$

$$P(MP = \text{Yes} | G = F) = \frac{3}{4}$$

$$P(WP = \text{Yes} | G = F) = \frac{2}{4}$$

$$P(CL = \text{No} | G = F) = \frac{1}{4}$$

$$P(CC = \text{No} | G = F) = \frac{3}{4}$$

$$\Rightarrow \frac{3}{4} \times \frac{2}{4} \times \frac{1}{4} \times \frac{3}{4} \times \frac{4}{10} = 0.281$$

Male has higher probability, Gender = "Male"

	Outlook	temperature	humidity	windy	play class.
S	sunny	88	85	85	NO
S	sunny	81	80	81	NO
M	overcast	83	83	86	YES
R	rainy	80	80	86	YES
R	rainy	68	80	80	YES
M	overcast	65	83	70	NO
M	sunny	64	65	80	YES
M	sunny	72	80	80	NO
M	rainy	69	80	70	YES
R	sunny	75	83	80	YES
R	overcast	72	90	70	YES
O	overcast	81	75	80	YES
O	rainy	71	91	70	NO

Find the values of the tuple, where attribute values
are as follows Output look = overcast, temperature = 16,
humidity = 62 windy = false.

write on desk

$$f(x) = \frac{1}{\alpha \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

R

$\mu - \text{temp} - \text{No} = 74.60$
 $\mu - \text{Hum} - \text{No} = 8.0$
 $\alpha - \text{temp} - \text{Yes} = 8.0$
 $\alpha - \text{hum} - \text{Yes} = 9.7$

$f(\text{temp} = 60 | \text{Yes}) = \frac{1}{6.2 \sqrt{2\pi}} e^{-\frac{(60-73)^2}{2(6.2)^2}}$
 $f(\text{temp} = 60 | \text{No}) = \frac{1}{8 \sqrt{2\pi}} e^{-\frac{(60-74.6)^2}{2(8)^2}}$
 $f(\text{hum} = 62 | \text{Yes}) = 0.0096$
 $f(\text{hum} = 62 | \text{No}) = 0.0018$

$P(\text{Yes}) = \frac{9}{14}$ $P(\text{No}) = \frac{5}{14}$
 $P(\text{Outlook} = \text{overcast} | \text{Yes}) = \frac{4}{9}$
 $P(\text{Outlook} = \text{overcast} | \text{No}) = \frac{0}{5}$
 $P(\text{Windy} = \text{false} | \text{Yes}) = \frac{6}{9}$
 $P(\text{Windy} = \text{false} | \text{No}) = \frac{2}{5}$

$\left. \begin{array}{l} \text{Poisson} \\ \text{Laplace} \end{array} \right\} \text{Laplace correction.}$

↳ "clustering" is the process of grouping the data into classes of cluster, so that objects within a cluster have high similarity in comparison to another, but are varying dissimilar to objects in other clusters.

Measurement of Similarity:

* categorical variable

* Oblique nominal variable.

→ Consider 2 instances x and y which consists of n nominal attributes then the distance b/w is measured as follows

$$d(x, y) = \frac{n-m}{n}$$

where $m = \text{no. of matches}$

→ consider 2 instances x and y that consists of n -numerical attributes. it can be represented as vector for $x_1, x_2, x_3, \dots, x_n$
 $y_1, y_2, y_3, \dots, y_n$

The distance b/w x and y will be determined

? by $d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$

Cosine-based similarity in the vector form - D.

$$\cos(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

Clustering methods:

→ partitioning methods

→ k-means

→ k-medoids / PAM. Partitioning Around medoids

k-means : It works only for numerical attributes

Squared error is used to determine clustering quality.

id	name	height	weight
x_1	Ram	64	60
x_2	Shyam	60	61
x_3	Gita	59	70
x_4	Mohan	68	71

$$d(x_3, x_1) = \sqrt{(59-64)^2 + (70-60)^2} = 11.48$$

$$d(x_4, x_1) = 11.4$$

$$d(x_4, x_2) = 12.81$$

$$E = d(x_2, x_3)^2 + d(x_4, x_1)^2 \approx 219$$

Iteration 2

$$C_1 \rightarrow \left(\frac{64+68}{2}, \frac{60+71}{2} \right) A \text{ and } C_2 \rightarrow \left(\frac{59+60}{2}, \frac{70+61}{2} \right) B$$

$$\text{Centroid for } C_1 = \text{mean } \frac{135}{2} = 67.5 \text{ (points)}$$

$$A = 66, 65.5$$

$$B = 59.5, 65.5$$

$$d(x_1, A) = 8.585$$

$$d(x_1, B) = 7.1$$

$$d(x_2, A) = 7.5$$

$$d(x_2, B) = 4.53$$

minimum distance from all points to cluster A

minimum distance from all points to cluster B

$$d(x_3, A) = 8.32$$

minimum distance from all points to cluster A

$$d(x_3, B) = 4.57$$

minimum distance from all points to cluster B

$$d(x_4, A) = 5.85$$

minimum distance from all points to cluster A

$$d(x_4, B) = 10.12$$

minimum distance from all points to cluster B

$$E = d(x_1, A)^2 + d(x_4, A)^2 + d(x_2, B)^2 + d(x_3, B)^2$$

$$= 109.49$$

minimum distance from all points to cluster A

minimum distance from all points to cluster B

⇒ This algo ensures that it selects initial k prototypes arbitrarily. The absolute error criteria is used to determine the clustering quality.

⇒ In each iteration the prototype of each cluster is assigned to an actual dataset of each point that minimizes the absolute error criteria.

Iteration 1:

$$E = d(x_2, x_3) + d(x_1, x_4) = 20.72$$

Iteration 2: The Algo. selects new midoids in place of existing midoids, i.e., either x_3 or x_4 selected in place of x_1 or x_2 .

⇒ consider the case where x_4 is selected in

place x_1

x_2

x_4

$$d(x_1, x_2) *$$

$$d(x_1, x_4) *$$

$$d(x_3, x_2)$$

$$d(x_3, x_4) *$$

$$\Sigma = d(x_1, x_2) + d(x_3, x_4).$$

$$14.8 \quad 7.8 \quad 0 \quad 7.8$$

$$0 \quad 7.8$$

$$0 \quad 14.8$$

$$0 \quad 7.8$$

$$0 \quad 14.8$$

$$0 \quad 7.8$$

$$0 \quad 0.1$$

$$0 \quad 7.8$$

$$0 \quad 14.8$$

$$0 \quad 7.8$$

- In this the first step to create similarity matrix is the pair of table of all pair wise distances or degree of similarity between clusters. Initially it contains pair wise distances between individual pair of records.
- Distance between the cluster can be found using 3 different approaches

1) Single linkage

2) complete linkage

3) Centroid linkage

