

Data Science

DS is an interdisciplinary field of scientific methods, processes and systems to extract unseen patterns, derive meaningful information and make business prediction. It uses complex ML algorithms to build predictive model.

Data Science Life Cycle

Step 1: Problem definition — ask relevant questions and define obj for the problem that needs to be tackled.

Step 2: Collecting data (data preparation & data investigation) — gather & scrap the raw data necessary for project formulation.

Step 3: Preprocessing data (Data-cleaning) — fix the inconsistencies within the data set & handle the missing values.

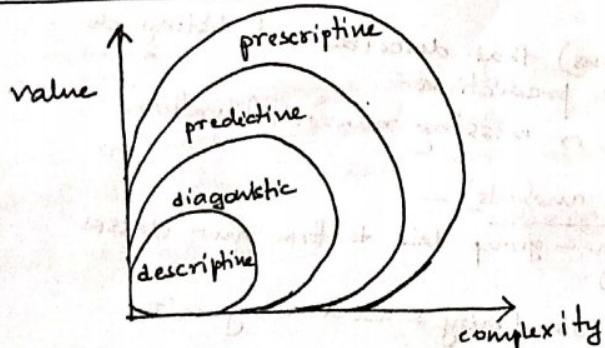
Step 4: Data exploration (EDA - Exploratory Data Analysis) MCDM -
Multi-criteria
Decision Making
Form hypothesis about your defined problem by visually analysing the data.

Step 5: Feature engineering — Select important feature and construct more meaningful ones using the data that we have. (PCA - Principle Component Analysis)

Step 6: Predictive modelling — train ML models, evaluate their performance, and use them to make predictions.

Step 7: Data visualization — communicate the findings using plots and interactive visualizations (eg. bar plot, pie plot, etc.)

Types of Data Analytics



• Descriptive analytics

It looks at data to examine, understand, and describe something that has already happened. (descriptive statistical analysis)

• Diagnostic analytics

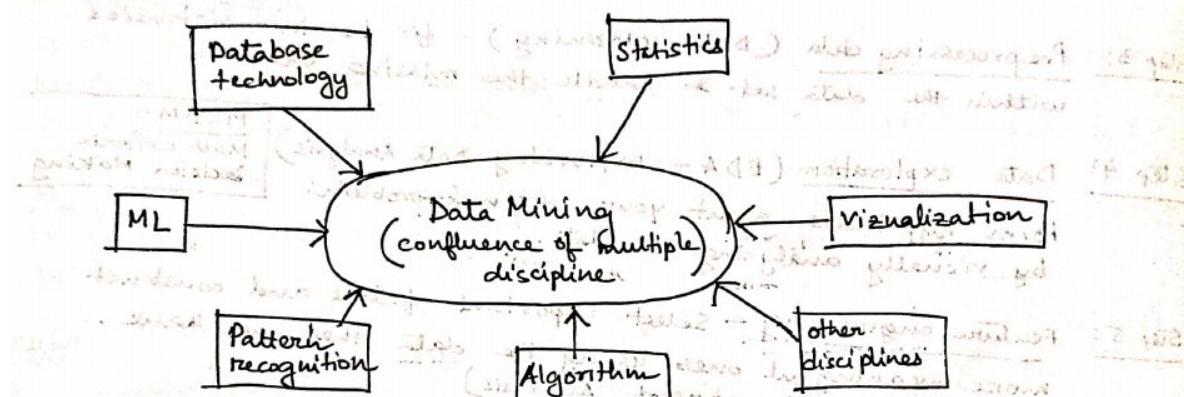
It goes deeper than descriptive analytics by seeking to understand the why behind what happened.

- **predictive analytics** - It relies on historical data, past trends and assumptions to answer questions about what will happen in the future.
- **prescriptive analytics** - It aims to identify specific actions that an individual/organisation should take to reach future targets or goals.

Data Mining

Knowledge Discovery from Databases (KDD) I

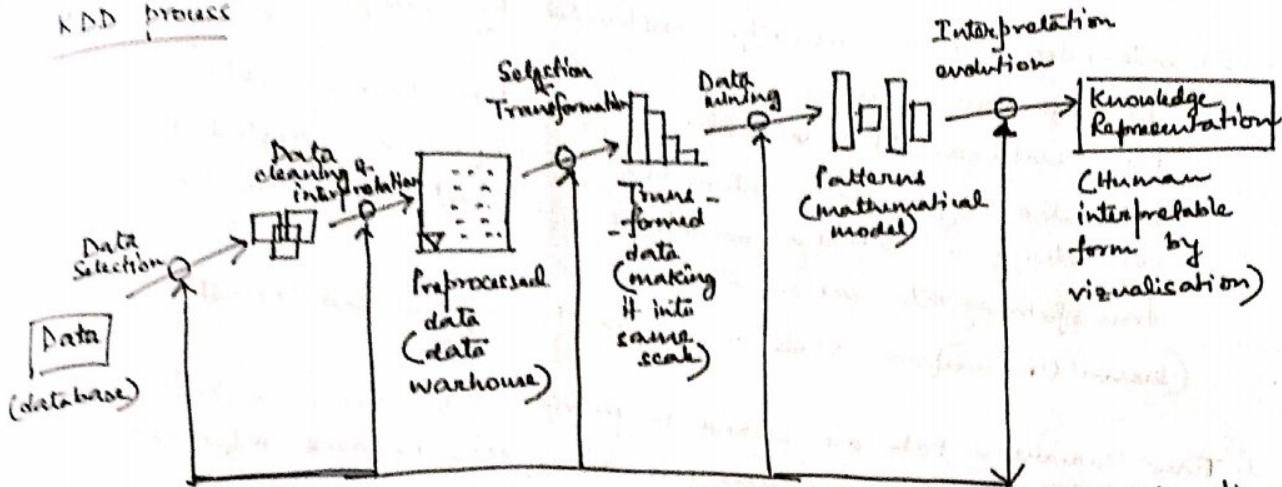
Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns/knowledge from huge amount of data.



Data Mining Functionalities

- Multi-dimensional concept description & characterization & discrimination.
Generalise, summarise and contrast data characteristics.
- frequent patterns, association, correlation (Association Rule Mining)
eg - market basket analysis
- classification & prediction -
 - (i) construct models (functions) that describe and distinguish classes/concepts for future prediction.
 - (ii) predict some unknown or missing numerical values.
- cluster analysis (exploratory analysis) -
 - (i) class label is unknown - group data to form new classes (unsupervised learning)
 - (ii) maximising intra class similarity & minimising inter class similarity.
- outlier analysis (anomaly detection) -
data that does not comply with the general behaviour of the dataset.
- trend & evolution analysis -
 - (i) trend & deviation
 - (ii) sequential pattern mining
 - (iii) periodicity analysis
 - (iv) similarity based analysis.
- other pattern directed / statistical analysis

KDD process



Knowledge discovery process as depicted in figure, consists of iterative sequence of steps.

1. Data Selection - Data relevant for analysis task is retrieved from db.

2. Data integration - multiple data sources may be combined.

3. Data cleaning - to remove noise & inconsistent data

4. Data transformation - where data are transformed or consolidated into forms appropriate for mining by performing summary or aggregation of.

5. Data mining - An essential process where intelligent methods are applied in order to extract data patterns.

6. Pattern evaluation - To identify the truly interesting patterns representing knowledge based on some interesting measures.

7. Knowledge representation - where visualisation & knowledge representation techniques are used to present the mined knowledge to the user.

Data Warehouse

It is the act of organising & storing data in a way so as to make its retrieval efficient and insightful. It is also called the process of transforming data into information.

According to Bill Inmon - Data warehouse is a

- subject oriented
- integrated
- non volatile
- time variant

collection of data in supports of management's decision making.

1. Subject Oriented - A data warehouse is organised around major subjects. It provides a simple & concise view around particular subject issues by excluding data that is not useful in the decision support process.

2. integrated - It is usually constructed by integrating multiple heterogeneous sources.

Eg - Relational DBs, online transaction records, etc.

Data cleaning & data integration techniques are applied to ensure consistency in naming conventions, attribute measures & also transforming into same scale.

(~~biased~~ Un-uniform scale results in having biased results)

3. Time Variant - Data are stored to provide info from a historical perspective. It allows analysis of past, related information to present and enables forecast to future.

4. Non-volatile - Data is not updated or deleted from warehouse in real time.

Difference b/w OLTP & OLAP

Features	Online Transaction Processing (OLTP)	Online Analytical Processing (OLAP)
characteristics	Operational processing	Informational processing
Orientation	Transaction analysis	Analysis of data
User	DB professionals	Knowledge worker (Data scientist, data analyst)
functions	Day to day operations	Long term informational requirements (Decision Support System) [DSS]
data	current transaction	Historical Data (Accuracy maintained over time)
Unit of work	Simple transaction queries	Complex queries
response time	requires faster response time.	Can have longer response time
access	read & write	mostly read
focus	data i/p	information out
No. of records accessed	100s of data	millions of data (big data)
No. of users	Thousands	Hundreds (specific requirement)
DB size	GB to higher order GB	> TB.

Functions of Data Warehouse Tools

(ETL - Extraction, Transformation, Loading)

- ① Data Extraction - gathers data from multiple external heterogeneous sources
- ② Data Transformation - converts data from host format to warehouse format
- ③ Data Loading - It sorts, summarizes, consolidates, computes, views, checks integrity and build indices & partitions.
- Besides E,T,L, 2 other terms are - (i) Data cleaning & (ii) data refreshing
They are important steps in improving the data quality & subsequent data mining results.

Association Rule Mining (ARM)

Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction.

Eg - Market Basket transaction

TID	Items	* unsupervised
1	Bread, Milk (k=2)	
2	Bread, Egg, Butter, Ketchup (k=4)	
3	Milk, Eggs, Butter, Cookies	
4	Bread, Milk, Eggs, Butter	
5	Bread, Milk, Eggs, Cookies	

Eg's of Association Rules

$$\{Eggs\} \Rightarrow \{Butter\}$$

$$\{Milk, Bread\} \Rightarrow \{Eggs\}$$

Antecedent
(LHS of AR)

Consequent
(RHS of AR)

⇒ Item Set^(k) - A collection of 1 or more items.

Rule Evaluation Metrics

1. Support count/frequency, count of the item set (σ) -

freq. of occurrence of an item set.

$\sigma \{Milk, Bread, Eggs\} = 2$ (Milk, Bread, Eggs occurring together in DS; TID=4 & 5)

2. Support / Relative Support^(s) - Fraction of transactions that contain an item set.

$$s \{Milk, Bread, Eggs\} = \frac{2}{5} = 0.40$$

↓ (Total no. of transactions in DS)

3. Frequent Item set - An item set, whose support is greater than or equals to a minimum support threshold (minsup threshold)

4. Confidence - It measures how often items in Y appear in transactions that contain X.

$$P(Y|X) = \text{confidence}(X \Rightarrow Y)$$

Range = [0, 1]

$$= \frac{\text{support}(X \cup Y)}{\text{support}(X)}$$

$\{ \text{Milk, Bread} \} \Rightarrow \{ \text{Eggs} \}$

$$\text{confidence} = \frac{\text{support}(\{ \text{Milk, Bread, Eggs} \})}{\text{support}(\{ \text{Milk, Bread} \})}$$

$$= \frac{0.4}{0.6} = \frac{2}{3} = 0.667$$

$\therefore 67\% \text{ of customers who purchased } \{ \text{Milk, Bread} \}, \text{ also purchased Eggs.}$

ARM Task

Given a set of transactions, T , the goal of ARM is to find all the rules having

- support \geq min sup threshold
- confidence \geq min conf threshold (min. confidence threshold)

Note: A major challenge in mining frequent item sets from a large dataset is the fact that such mining often generates a huge no. of item sets satisfying the minimum support threshold, especially when min. support is set low. This is because, if an item set is frequent, each of its subsets is frequent as well. A long item set will contain a combinatorial no. of shorter, frequent sub item sets.

Lift (Correlation Measure) -

Support & confidence measures are insufficient at filtering out uninteresting association rules. To tackle this weakness, a correlation measure can be used to augment the support-confidence framework for association rules.

$$\text{lift}(A, B) = \frac{\text{Support}(A \cup B)}{\text{Support}(A) \cdot \text{Support}(B)}$$

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

There are 3 cases :-

- (i) lift < 1 (occurrence of A is negatively co-related with that of B)
- (ii) lift > 1 (+ve co-relation)
- (iii) lift $= 1$ (confidence = support)
(A & B are independent and there is no co-relation b/w them)

TDB

- Q. A TDB has 500 transactions out of which X appears in 77 transactions and out of these 77 transactions another item 'Y' appears in 35 transactions. Find the confidence of the rule $X \rightarrow Y$

$$\text{confidence} = \frac{\text{Support}(X \cup Y)}{\text{Support}(X)} = \frac{\frac{35}{500}}{\frac{77}{500}} = \frac{35}{77} = 0.454$$

- Q. let $\text{support}(X) = 0.2$, & $\text{support}(Y) = 0.4$ and $\text{support}(X \cup Y) = 0.1$
Find lift (X, Y) ?

$$\therefore \text{lift}(X, Y) = \frac{\text{support}(X \cup Y)}{\text{support}(X) \text{support}(Y)} = \frac{0.1}{0.2 \times 0.4} = \frac{0.1}{0.08} = 1.25$$

- Q. A TDB has 6 items, such as U, V, W, X, Y, Z. Out of 10 transactions U, V, W, X, Y, Z appears in 2, 5, 7, 3, 4 and 6 transactions respectively. Let the min^{um} support threshold is 0.5. Find the list of rare items.
 $\text{Support}(U) = 0.2$, $\text{sup}(V) = 0.5$, $\text{sup}(W) = 0.7$, $\text{sup}(X) = 0.3$, $\text{sup}(Y) = 0.4$, $\text{sup}(Z) = 0.6$
 $\therefore U, X, Y$ are rare items.

Apriori Algorithm

- Apriori Principle — If an item set is frequent, then all of its subsets must also be frequent (Reducing no. of candidate item sets)
Eg- If $\{A, B, C, D\}$ is a freq. item set, then any subset like $\{A, B, C\}$ or $\{B, D\}$ is also a freq item set.
Apriori Principle holds due to the following property of the support measure:
- $\forall X, Y : (X \subseteq Y) \Rightarrow \text{Support}(X) \geq \text{Support}(Y)$
- Anti monotone property of support.
Support of an item set never exceeds the support of its subsets.

Apriori Algo

1. let $k=1$
2. Generate frequent itemsets of length 1 (1-itemset)
3. Repeat until no frequent itemsets are identified
 - i. Generate length $(k+1)$ candidate itemsets from length k frequent itemsets
 - ii. Prune candidate itemsets containing subsets of length k that are infrequent
 - iii. Count the support of each candidate by scanning the database.
 - iv. Eliminate candidates that are infrequent leaving only those that are frequent.
4. Generate strong rules (based on confidence α value)

Q. For the following given transaction dataset, generate rules using Apriori Algo.

Consider the values of support = 22% and confidence = 70%.

Transaction ID | Items purchased

1.	I_1, I_2, I_5	(I_1, I_2) happens \Leftrightarrow (I_1, I_5)
2.	I_2, I_4	(I_2, I_4) happens \Leftrightarrow (I_2, I_3)
3.	I_2, I_3	(I_2, I_3) happens \Leftrightarrow (I_1, I_3)
4.	I_1, I_2, I_4	(I_1, I_2) happens \Leftrightarrow (I_1, I_3)
5.	I_1, I_3	(I_1, I_3) happens \Leftrightarrow (I_1, I_2)
6.	I_2, I_3	(I_2, I_3) happens \Leftrightarrow (I_1, I_3)
7.	I_1, I_3	(I_1, I_3) happens \Leftrightarrow (I_1, I_2)
8.	I_1, I_2, I_3, I_5	(I_1, I_2, I_3) happens \Leftrightarrow (I_1, I_2, I_5)
9.	I_1, I_2, I_3	(I_1, I_2, I_3) happens \Leftrightarrow (I_1, I_2, I_5)

Sol:

Step 1 1 frequent candidate item sets ($k=1$) or (C_1)

Item	frequency	Relative Support %.
I_1	6	$6/9 = 66.67\%$
I_2	7	$7/9 = 77.78\%$
I_3	6	$6/9 = 66.67\%$
I_4	2	$2/9 = 22.22\%$
I_5	2	$2/9 = 22.22\%$

all the items have support $\geq 22\%$.

Step 2 Generate pairs of item sets
2 candidate item sets ($k=2$) or (C_2)

Itemset pairs	Freq	Relative support
$\{I_1, I_2\}$	4	44.44%
$\{I_1, I_3\}$	4	44.44%
$\times \{I_1, I_4\}$	1	11.11%
$\{I_1, I_5\}$	2	22.22%
$\{I_2, I_3\}$	4	44.44%
$\{I_2, I_4\}$	2	22.22%
$\{I_2, I_5\}$	2	22.22%
$\times \{I_3, I_4\}$	0	0%
$\times \{I_3, I_5\}$	1	11.11%
$\times \{I_4, I_5\}$	0	0%



Step 3 3 candidate item sets (k=3) or C_3

Itemset	frequency	Relative support
$\{I_1, I_2, I_3\}$	2	22.23%
$\{I_1, I_2, I_5\}$	2	22.23%

Step 4 4 candidate item sets (k=4) or C_4

Item set	frequency	relative support
$\{I_1, I_2, I_3, I_5\}$	1	11.11%
$\{I_1, I_2, I_3, I_6\}$	1	11.11%

\therefore Algo terminates
freq. item sets = $\{I_1, I_2, I_3\}$, $\{I_1, I_2, I_5\}$
obtained:

Step 5 Generate the rules based on the frequent item sets

(i) rules generated by $\{I_1, I_2, I_5\}$

$$\begin{aligned}
 (I_1, I_2) \rightarrow I_5 &= \frac{22.23}{44.45} = 50\% \times \\
 (I_1, I_5) \rightarrow I_2 &= \frac{22.23}{22.23} = 100\% \times \\
 (I_2, I_5) \rightarrow I_1 &= \frac{22.23}{22.23} = 100\% \times \\
 I_1 \rightarrow (I_2, I_5) &= \frac{22.23}{66.67} = 33.33\% = 33\% \times \\
 I_2 \rightarrow (I_1, I_5) &= \frac{22.23}{77.78} = 28.57 = 29\% \times \\
 I_5 \rightarrow (I_1, I_2) &= \frac{22.23}{22.23} = 100\% \times
 \end{aligned}$$

rules generated:

$$\begin{aligned}
 (I_1, I_5) \rightarrow I_2 \\
 (I_2, I_5) \rightarrow I_1 \\
 I_5 \rightarrow (I_1, I_2)
 \end{aligned}$$

(ii) rules generated by $\{I_1, I_2, I_3\}$

$$\begin{aligned}
 I_1 \rightarrow (I_2, I_3) &= \frac{2/6}{6} = 33\% \\
 I_2 \rightarrow (I_1, I_3) &= \frac{2/7}{7} = 29\% \\
 I_3 \rightarrow (I_1, I_2) &= \frac{2/6}{6} = 33\% \\
 (I_1, I_2) \rightarrow I_3 &= \frac{2/4}{4} = 50\% \\
 (I_1, I_3) \rightarrow I_2 &= \frac{2/4}{4} = 50\% \\
 (I_2, I_3) \rightarrow I_1 &= \frac{2/4}{4} = 50\%
 \end{aligned}$$

for $\{I_1, I_2, I_3\}$ no rules are getting generated

\therefore for calculating the support of itemset $\{I_1, I_2, I_3\}$ we have to calculate the support of itemset $\{I_1, I_2, I_5\}$ and $\{I_1, I_2, I_6\}$ which is 11.11%.

FP-Growth Algorithm

Q. Can we design any method that mines the complete set of freq item sets without candidate generation?

Sol: An interesting method in this attempt is called FP-Growth (frequent-pattern growth)

- It adapts a divide & conquer strategy as follows:
 - Firstly, it compresses the database representing frequent items into a frequent pattern tree which retains the itemset association information.
 - It then divides the compressed database into a set of conditional databases (projected database), each associated with one freq. item and mines each search database separately.

Q1. Generate FP-Tree for the following transaction dataset whose minimum support count = 3. Show the conditional pattern base, conditional FP-tree and frequent item set.

Transaction ID	Items
T ₁	{E, K, M, N, O, Y}
T ₂	{D, E, K, N, O, Y}
T ₃	{A, E, K, M}
T ₄	{C, K, M, U, Y}
T ₅	{C, E, I, K, O}

Step 1

The first step is to scan the database to find the occurrences of the item sets in the database. This step is same as the first step of Apriori Algo.

Item	frequency
A	1 X
C	2 X
D	1 X
E	4
I	1 X
K	5
M	3
N	2 X
O	3
U	1 X
Y	3

Step 2
A frequent pattern set (L) is built which will contain all the elements whose frequency is greater than or equal to the minimum support. These elements are stored in descending order of their respective frequencies.

$$L = \{ K:5, E:4, M:3, O:3, Y:3 \}$$

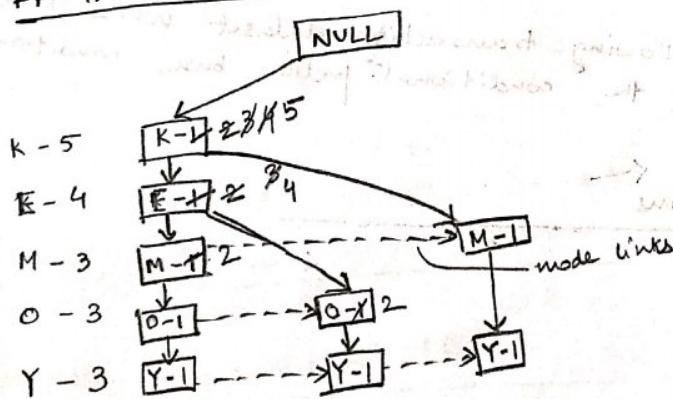
appear
in
data set first

Step 3 Now for each transaction, the respective ordered item set is built:

Transaction ID	Ordered item set (descending order)
T ₁	{K, E, M, O, Y}
T ₂	{K, E, O, Y}
T ₃	{K, E, M}
T ₄	{K, M, Y}
T ₅	{K, E, O}

Step 4: Now all the ordered item sets are inserted into a Trie Data structure. To facilitate tree traversal, an item header table is built so that, each item points to its occurrences in the tree via a chain of node links. In this way the problem of mining frequent patterns in db is transformed to that of mining the fp-tree.

FP-Tree Construction



Step 5: Next step is to mine the created FP-tree. Now for each item, the conditional pattern base is computed which is path ~~leads~~ ^{tags} labels of all the paths which ~~leads~~ ^{leads} to any node of the given item in the FP tree.

Items	Conditional pattern base
Y	{K, E, M, O:1}, {K, E, O:1}, {K, M:1}
O	{K, E, M:1}, {K, E:2}
M	{K, E:2}, {K:1}
E	{K:4}
K	-

Step 6

Now for each item, the 'conditional FP tree is built.

Items	Conditional FP Tree
Y	$\{K:3\}$ ← check frequency & support count & find the items which satisfy it
O	$\{K, E:3\}$
M	$\{K:3\}$
E	$\{K:4\}$
K	—

Step 7

from the conditional FP Tree, the freq. pattern rules are generated by pairing the items of the conditional freq. pattern tree sets to the corresponding item as shown below:

Items	Frequent pattern generated
Y	$\{\langle K, Y:3 \rangle\}$
O	$\{\langle K, O:3 \rangle, \langle E, O:3 \rangle, \langle K, E, O:3 \rangle\}$
M	$\{\langle K, M:3 \rangle\}$
E	$\{\langle K, E:4 \rangle\}$
K	—

Q. Generate FP Tree for the following transaction dataset whose minimum support count is 2. Show the conditional pattern base, conditional FP tree & freq. item set.

TID	List of items
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

Sol.

S1	Item	Frequency
I1	6	
I2	7	
I3	6	
I4	2	
I5	2	

S2

$$L = \{I2:7, I1:6, I3:6, I5:2, I4:2\}$$

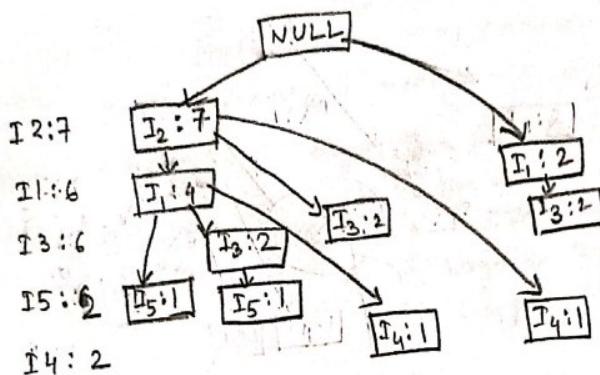


Scanned with OKEN Scanner

53

Transaction id	Ordered item set
T100	{I ₂ , I ₁ , I ₅ }
T200	{I ₂ , I ₄ }
T300	{I ₂ , I ₃ }
T400	{I ₂ , I ₁ , I ₄ }
T500	{I ₁ , I ₃ }
T600	{I ₂ , I ₃ }
T700	{I ₁ , I ₃ }
T800	{I ₂ , I ₁ , I ₃ , I ₅ }
T900	{I ₂ , I ₁ , I ₃ }

54



55

Items	conditional pattern base
I ₄	{I ₂ , I ₁ :1}, {I ₂ :1}
I ₅	{I ₂ , I ₁ :1}, {I ₂ , I ₁ , I ₃ :1}
I ₃	{I ₂ :2}, {I ₁ :2}, {I ₂ , I ₁ :2}
I ₁	{I ₂ :4}, left subtree right subtree
I ₂	—

56

Items	conditional fp Tree
I ₄	{I ₂ :2}
I ₅	{I ₂ :2}, {I ₁ :2}
I ₃	{I ₂ :4}, {I ₁ :2}, {I ₁ :2} (treat subtree diff) left subtree right subtree
I ₁	{I ₂ :4}
I ₂	—

57

Items	frequent pattern generated
I ₄	{<I ₂ , I ₄ :2>}
I ₅	{<I ₂ , I ₅ :2>}, <I ₁ , I ₅ :2>, <I ₂ , I ₁ , I ₅ :2>}
I ₃	{<I ₂ , I ₃ :4>}, <I ₁ , I ₃ :4>, <I ₁ , I ₂ , I ₃ :2>}
I ₁	{<I ₁ , I ₂ :4>}
I ₂	—

min of {I₂:4,
I₁:2} &
{I₁:2}

Adv:-

- Algo needs to scan the database twice when compared to Apriori which scans the ~~first~~ transactions for each iterations.
- No candidate item generation, hence algo takes less amount of time.
- Requires less memory in comparison to Apriori (for large db specially)
- Efficient & scalable for mining both long & short freq. patterns

Disadv:-

- building of FP tree is more complex than Apriori
- This algo ^{is} ~~may be~~ expensive for small databases

CLASSIFICATION AND PREDICTION (Module-2)

- Classification & Prediction are two forms of data analysis that can be used to extract models that describe important data class or predict future data trends. Such analysis can help to provide us with a better understanding of data at large.
- Classification predicts categorical (discrete) labels, prediction models continuous valued func.

Ex. We can build a classification model to categorise bank loan applications as either safe / risky, & a prediction model predicts the expenditure of customers given their income & occupation.

Popular classification algo: ① K - nearest Neighbours (KNN)
② Decision Tree
③ Naive Bayes classifier
④ Logistic Regression
⑤ Support Vector Machine

Supervised & Unsupervised Learning

Supervised Learning

- ① These algo are trained using labelled data.
- ② The goal is to train the model so that it can predict the o/p when it is given new data.
- ③ It can be used for those cases where we know i/p & o/p.
- ④ It can be categorised in classification & regression problems.
- ⑤ Computational complexity less complex than unsupervised learning.
- ⑥ These models produce more accurate result generally.
- ⑦ Linear Regression, Logistic regression, Bayesian logic, etc.

Unsupervised Learning

- ① These algo are trained using unlabelled data.
- ② The goal is to find the hidden patterns and useful insights from the unknown dataset.
- ③ It can be used for those cases where we have i/p data only.
- ④ Classified in clustering & association problems.
- ⑤ Higher computational complexity
- ⑥ These models give less accurate result as compared to supervised learning.
- ⑦ Clustering, Apriori algo, etc.

Cluster Analysis

- It is the process of grouping the data into classes / clusters so that objects within a cluster have high similarity in comparison to one another but are very dissimilar to obj in other clusters.
- Cluster analysis have been widely used in numerous applications like market research, pattern recognition, image processing, etc.
- It is also called data segmentation becuz clustering partitions large data sets into groups acc. to their similarity.
- Clustering can also be used for outlier detection where outliers (values that are far away from any cluster) may be more interesting than common ones.

- It is an example of unsupervised learning.

[learning by examples - supervised
" " observation - unsupervised]

Categorisation of major clustering methods

* (i) Partitioning methods :

Given a database of n objects, a partitioning method constructs k -partitions of the data, where each partition represents a cluster and $k \leq n$. It classifies the data into k -groups, which satisfies the following requirements:

1. Each group must contain atleast one object.
2. Each object must belong to exactly one group.

* Centroid Based Technique - K-Means -

1. K-means algo takes the i/p parameter k and partitions a set of n -objects into k -clusters so that the resulting intra cluster similarity is high for the inter-cluster similarity is low.
2. cluster similarity is measured in regard to the mean value of the objects in a cluster which is referred to as cluster's centroid.

K-Means clustering algorithm

Input: $k \rightarrow$ no. of clusters
 $D \rightarrow$ dataset containing n objects $[x_1, x_2, \dots, x_n]$

Output: A set of k clusters.

Step 1 - choose k random points as initial cluster centers called centroids.

Step 2 - Assign each x_i to the closest cluster by implementing its distance to each centroid.

Step 3 - for each cluster, new centroids are computed by taking the avg of the assigned points. (updating the cluster mean)

Step 4 - Keep repeating step 2 & 3 until convergence is achieved.

* Note: - K-means clustering stops creating/optimizing clusters after finding no new reassignment of data points and also after the algo reaches the defined no. of iterations.

- * Q1. Apply the K-means algo for the following dataset for 2 clusters. Consider data points 1 and 2 as the initial centroid of the respective clusters. Find out the final cluster centroids. Continue the procedure for 2 iterations. Apply Euclidean dist. for as distance function.

Sample no.	X	Y
1	185	72
2	170	56
3	168	60
4	179	68
5	182	72
6	188	77

10-15 marks \rightarrow maxⁿ Iteration 3
 otherwise \rightarrow " " 2

Point	Centroid 1 (185, 72)	cluster 2 (170, 56)	cluster assignment
(185, 72)	0	21.93	c_1
(170, 56)	21.93	0	c_2
(168, 60)	20.81	4.47	c_2
(179, 68)	7.21	15	c_1
(182, 72)	3	20	c_1
(188, 77)	5.83	27.65	c_1

New centroid

$$c_1: \left(\frac{185 + 179 + 182 + 188}{4}, \frac{72 + 68 + 72 + 77}{4} \right) \\ = (183.5, 72.25)$$

$$c_2: \left(\frac{170 + 168}{2}, \frac{56 + 60}{2} \right) \\ = (169, 58)$$

Iteration-II

After the 2nd iteration, the assignment of points has not changed & hence the algorithm is stopped & the points are clustered.

final centroid value: $c_1 = 1, 4, 5, 6$
 $c_2 = 2, 3$

	P_1	P_2	P_3	P_4	P_5
P_1	0				
P_2	9	0			
P_3	3	7	0		
P_4	6	5	9	0	
P_5	11	10	(2)	8	0
			min		

	P_1	P_2	$[P_3, P_5]$	P_4
P_1	0			
P_2	9	0		
P_3	11	10	0	
P_4	6	(5) min	9	0

$\downarrow [P_1, P_2, P_4], [P_3, P_5]$

Pair $(P_1, (P_2, P_4))$
dist = 9

	P_1	$[P_2, P_4]$	$[P_3, P_5]$
P_1	0		
$[P_2, P_4]$	9 min	0	

$\rightarrow \cancel{d((P_1, P_2), (P_3, P_5))}, \cancel{d((P_1, P_3), (P_2, P_5))},$

$$\begin{aligned} & \max(d((P_3, P_5), P_2), \\ & \quad d((P_3, P_5), P_4)) \\ & = \max(10, 9) \\ & = 9 \end{aligned}$$

Average Linkage

The hierarchical method that involves looking at the distances b/w all pairs and avg. all of these distances.

Q. for the given set of points, identify clusters using avg. linkage agglomerative clustering. Round off the dist. matrix to two places of decimal.

	A	B
$P_1.$	1	1
$P_2.$	1.5	1.5
$P_3.$	5	5
$P_4.$	3	4
$P_5.$	4	4
$P_6.$	3	3.5

	P_1	P_2	P_3	P_4	P_5	P_6
P_1	0					
P_2	0.71	0				
P_3	5.66	4.95	0			
P_4	3.60	2.92	2.24	0		
P_5	4.24	3.53	1.41	1.0	0	
P_6	3.20	2.5	2.5	0.5	1.12	0

Pair (P_4, P_6)
dist = 0.5

	P_1	P_2	P_3	P_4	P_5
P_1	0				
P_2	0.71	0			
P_3	5.66	4.95	0		
P_4	3.40	2.71	2.37	0	
P_5	4.24	3.53	1.41	1.06	0

↓ Pair (P_1, P_2)
↓ dist = 0.71

	$[P_1, P_2]$	P_3	$[P_4, P_5, P_6]$
$[P_1, P_2]$	0		
P_3	5.31	0	
$[P_4, P_5, P_6]$	3.47	1.89	0

Pair $(P_5, (P_4, P_6))$
dist = 1.06

	$[P_1, P_2]$	P_3	$[P_4, P_6]$	P_5
$[P_1, P_2]$	0			
P_3	5.31	0		
$[P_4, P_6]$	3.05	2.37	0	
P_5	3.89	1.41	1.06	0

	$[P_1, P_2]$	$[P_3, P_4, P_5, P_6]$
$[P_1, P_2]$	0	
$[P_3, P_4, P_5, P_6]$	4.39	0

4 Axioms of Distance Metrics

- (i) $d(i, j) \geq 0$ (distance is a non-negative no.)
- (ii) $d(i, i) = 0$ (distance of an obj. to itself is always 0)
- (iii) $d(i, j) = d(j, i)$ (distance is a symmetric func.)
- (iv) $d(i, j) \leq d(i, k) + d(k, j)$ (triangle inequality)
[For 3 elements in this set, the sum of the dist for any two pairs must be greater than the dist for the remaining part]

* * weighted Euclidean dist.

$$d(i, j) = \sqrt{w_1(x_{i1} - x_{j1})^2 + w_2(x_{i2} - x_{j2})^2 + \dots + w_n(x_{in} - x_{jn})^2}$$

K-Medoids Clustering

- (i) K-means algo is sensitive to outliers because an object with an extremely large value may substantially distort the distribution of data. Instead of taking the mean value of the obj. in the cluster as a reference point, we can pick actual objects to represent the clusters.
- (ii) Each remaining ~~object~~ object is clustered with the representative obj to which it is the most similar. The partitioning method is then performed based on the principle of minimising the sum of the dissimilarities b/w each object & its corresponding reference point.

$$E = \sum_{j=1}^k \sum_{p \in C_j} |p - o_j|$$

$E \rightarrow$ the sum of the absolute error for all objects in the dataset.

$p \rightarrow$ point in space representing a given object. in cluster C_j

$o_j \rightarrow$ Representative obj. of C_j (medoid)

Algo (PAM- Partitioning Around Medoids)

I/p: $k \rightarrow$ the no. of clusters

$D \rightarrow$ the dataset containing 'n' objects

O/p: A set of 'k' clusters.

Step 1: Initially select k -random points as the medoids from the given 'n' data points of the dataset.

Step 2: Associate each data point to the closest medoid.

Step 3: While the cost of the configuration decreases :

for each medoid 'm' and for each non-medoid datapoint 'o'

(i) swap 'm' and 'o', recompute the cost (sum of dist. of points to their medoid)

(ii) If the total cost of the configuration increase than the previous step, undo the swap.

Q. QB (42)

Step 1 Select 2 medoids: $c_1 = (3, 4)$, $c_2 = (7, 4)$

(max 3 iteration
min 2 iteration)

i	X	Y	c_1	c_2	cluster assignment
x_1	2	6	3	7	c_1
x_2	3	4	0	4	c_1
x_3	3	8	4	8	c_1
x_4	4	7	4	6	c_1
x_5	6	2	5	3	c_2
x_6	6	4	3	1	c_2
x_7	7	3	5	1	c_2
x_8	7	4	4	0	c_2
x_9	8	5	0	2	c_2
x_{10}	7	6	6	2	c_2

clusters are $\rightarrow c_1 = \{(2, 6), (3, 4), (3, 8), (4, 7)\}$
 $c_2 = \{(6, 2), (6, 4), (7, 3), (7, 4), (8, 5), (7, 6)\}$

Step 2
total cost

$$(3+4+4) + (3+1+1+2+2) = 20$$

$\begin{matrix} \text{dist} & \text{dist} \\ \text{b1w} & \text{b1w} \\ (3,4) & (3,8) \\ (2,6) & (4,7) \end{matrix}$

Step 3
Next step is the selection of non medoid o' randomly. Let us
assume $o' = (7, 3)$

$$\therefore c_1 = (3, 4)$$

$$o' = (7, 3)$$

i	X	Y	c_1	o'	
x_1	2	6	3	8	c_1
x_2	3	4	0	5	c_1
x_3	3	8	4	9	c_1
x_4	4	7	4	7	c_1
x_5	6	2	5	2	c_2
x_6	6	4	3	2	c_2
x_7	7	3	5	0	c_2
x_8	7	4	4	1	c_2
x_9	8	5	0	3	c_2
x_{10}	7	6	6	3	c_2

$$\begin{aligned} \text{Total cost} \\ = (3+4+4) + (2+2+1+3+3) \\ = 22 \end{aligned}$$

Step 4
Cost of swapping medoid c_2 with o'
 $s = \text{current total cost} - \text{previous total cost}$
 $\geq 22 - 20 \geq 2 > 0$

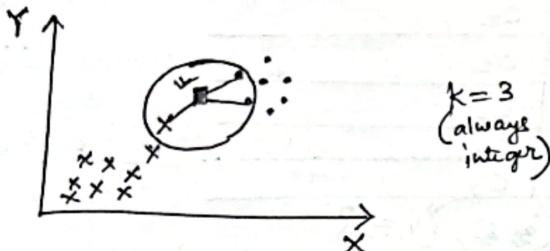
Moving to o' would
not be suitable one
as the prev one was
better
So final medoid = $(3, 4), (7, 3)$

K-NN Technique

(K-Nearest Neighbor)

[supervised algorithm]

- K-NN is a supervised classification algorithm that classifies the new data points based on the nearest data point.



{ for regression - use mean
 for classification - use mode

- K-NN algo is used to solve classification & regression problems. However it is mainly used for classification problem.
- K-NN is also called a lazy learning algo (lazy learner) because it doesn't perform any training when you supply the training data. Instead it just stores the data during training time and doesn't perform any calculations.
- K-NN is also called non-parametric method because it doesn't make any assumption about the underlying data distributions.

Algo

Step 1 - Load the dataset

Step 2 - choose the value of K (nearest data points)

Step 3 - for each point in the test dataset, do the following -

3.1 - calculate the dist b/w test data and each row of training data with the help of any dist. metric.

3.2 - Now based on the dist. value sort them in ascending order

3.3 - Next, it will choose the top K entries from the sorted list

3.4 - Label the test point based on the majority of classes present in the selected points.

Step 4 - Stop.

- ① Perform the K-NN classification algo on the following dataset and predict the class for x ($P_1 = 3, P_2 = 7$)

Given $K = 3$

	P_1	P_2	class
(i)	7	7	F
(ii)	7	4	F
(iii)	3	4	T
(iv)	1	4	T

Sol: $D(x, i) = \sqrt{(3-7)^2 + (7-7)^2} = 4 \rightarrow N_1$ False

$D(x, ii) = \sqrt{(3-7)^2 + (4-7)^2} = 5$

$D(x, iii) = \sqrt{(3-3)^2 + (4-7)^2} = 4 \rightarrow N_1$ True

$D(x, iv) = \sqrt{(3-1)^2 + (7-4)^2} = \sqrt{4+9} = \sqrt{13} = 3.61$

2 True > 1 False

N₂
True

$\therefore x (P_1=3, P_2=7)$ belongs to True class.

- Note:
- (i) when dealing with 2 class problem, it's better to choose an odd value of K .
 - (ii) otherwise a scenario can arise where the no. of neighbours in each class is the same.
 - (iii) choosing a very small value of K leads to unstable decision boundaries. For eg. - Suppose $K=1$ or 2 can be subjected to noise and effects of outliers]

Disadv:

- (i) associated computation cost is high as it stores all the training data
- (ii) prediction may become slower if the value of N is high.
- (iii) sensitive to irrelevant features.

Decision Tree Classification

- It is a supervised learning algorithm that can be used for both classification & regression. But mostly it is preferred for solving classification problems.
- It is a tree structured classifier where each internal node tests an attribute, each branch corresponds to attribute values and each leaf node assigns a classification.
- It is a graphical representation for getting all the possible solutions to a problem / decision based on given conditions.

Attribute Selection Measure (ASM)

While implementing a decision tree the main issue arises that how to select the best attribute of the root node. With the help of ASM technique, we can easily select the best attribute for the nodes of the tree.

Two popular technique : - (i) Entropy & Information Gain (ii) Gini index.

Entropy & Information Gain

(i) Entropy & Informational Gain

Entropy - It is a metric to measure the impurity in a given attribute. It specifies randomness in data.

$$\text{Entropy } (S) = \sum_{i=1}^n -P_i \log_2 (P_i)$$

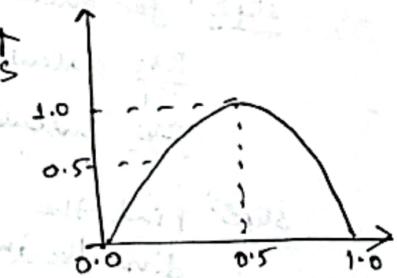
$$= -P_{\oplus} \log_2 (P_{\oplus}) - P_{\ominus} \log_2 (P_{\ominus})$$

proportion of the examples in S .

• Entropy value 0 represents data sample is

pure / homogeneous

• Entropy value 1 represents that the data sample has a 50-50 split belonging to 2 categories.



GPA	Studied	Passed
Low	F	F
Low	T	T
Medium	F	F
Medium	T	T
High	F	T
High	T	T

Q. what is the ~~ent~~ Entropy of the data set.

$$\text{Passed} = \begin{cases} F: 2 \rightarrow 4+, 2- \\ T: 4 \end{cases}$$

$$P_{\oplus} = \frac{4}{6}$$

$$\begin{aligned} \text{Entropy} &= -P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus} \\ &= -\frac{2}{6} \log_2 \left(\frac{2}{6}\right) - \frac{4}{6} \log_2 \left(\frac{4}{6}\right) \\ &= -\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) \\ &= 0.918 \approx 0.92 \end{aligned}$$

$$P_{\ominus} = \frac{2}{6}$$

Information Gain (IG) - It is the measurement of changes in entropy after the segmentation of a data set based on an attribute. It calculates how much information a feature provides to us about a class. According to the value of IG , we split the node and build the decision tree. A decision tree algorithm always tries to maximise the value of IG and a node or attribute having the highest IG is splitted first.

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{\text{values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

Q. for the following weather attributes create decision tree. Sample Dataset is given where the target attribute is "PlayTennis". $\circledast \rightarrow$

formation of decision Tree (steps in ID3)

Iterative Dicotomiser

Step 1: calculate the entropy of the entire D^S .

Step 2: for each attribute/feature:

2a: calculate entropy for all of its categorical values

2b: calculate information ~~of~~ gain for the feature.

Step 3: Find the feature with max^m information gain. Database is divided into smaller subsets acc. to categories of decision rules.

Step 4: Repeat it until we get the desired tree.

① →
Sol:

Entropy of entire dataset (for target attribute) → 9 Yes, 5 No
= [9+, 5-]

$$\text{entropy} = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right)$$
$$= 0.9403$$

Outlook → sunny, overcast, rain
(5) (4) (5)
[2+, 3-] [4+, 0-] [3+, 2-]

$$S(\text{outlook, sunny}) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) = 0.971$$

$$S(\text{outlook, overcast}) = 0$$

$$S(\text{outlook, rain}) = -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) = 0.971$$

$$\text{Information Gain}(\text{outlook}) = 0.940 - \left(\frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.971 \right)$$
$$= 0.940 - 0.694$$
$$= 0.246 \rightarrow \text{max}^m \rightarrow \text{root}$$

Temperature → ~~Hot~~ Hot, cool, Mild
(4) (4) (5)
[2+, 2-] [3+, 1-] [4+, 2-]

$$S(\text{Temp, hot}) = 1.0$$

$$S(\text{Temp, cool}) = 0.8113$$

$$S(\text{Temp, mild}) = 0.9183$$

$$IG(\text{Temp}) = 0.94 - \left(\frac{4}{14} \times 1 + \frac{4}{14} \times 0.8113 + \frac{6}{14} \times 0.9183 \right)$$
$$= 0.0289$$

Humidity → High, Normal
(7) (7)
[3+, 4-] [4+, 1-]

$$S(\text{Humidity, high}) = 0.9852$$

$$S(\text{ " , normal}) = 0.5916$$

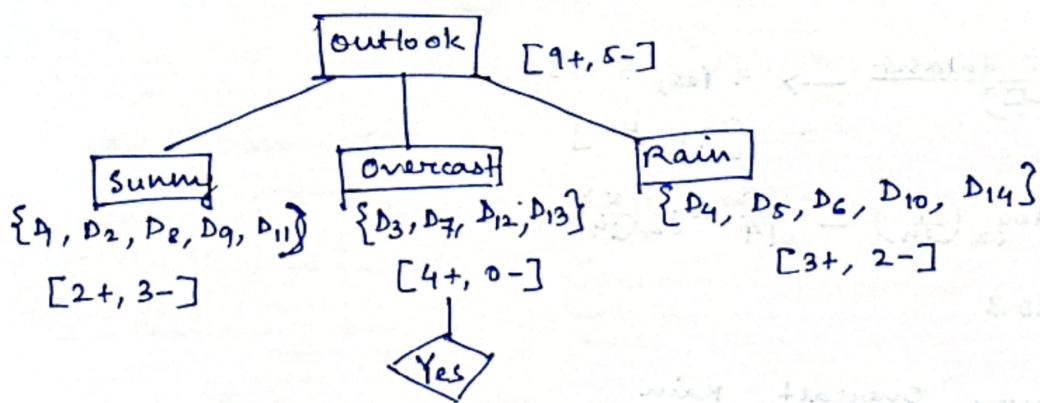
$$IG(\text{humidity}) = 0.1516$$

Wind → strong, weak
(6) (8)
[3+, 3-] [4+, 2-]

$$S(\text{wind, strong}) = 1$$

$$S(\text{wind, weak}) = 0.8113$$

$$IG(\text{wind}) = 0.0478$$



Partial Decision Tree

for Sunny

$$S(\text{Sunny}) [2+, 3-] = 0.971$$

temperature

$$\rightarrow \text{Shot} \rightarrow [0+, 2-] = 0$$

$$\rightarrow \text{Mild} \rightarrow [1+, 1-] = 1$$

$$\rightarrow \text{Cold} \rightarrow [1+, 0-] = 0$$

$$IG(\text{Sunny, temperature}) = 0.971 - \left(\frac{3}{5} \times 0 + \frac{2}{5} \times 1 + \frac{1}{5} \times 0 \right) = 0.571$$

humidity

$$\rightarrow \text{High} \rightarrow [0+, 3-] = 0$$

$$\rightarrow \text{Normal} \rightarrow [2+, 0-] = 0$$

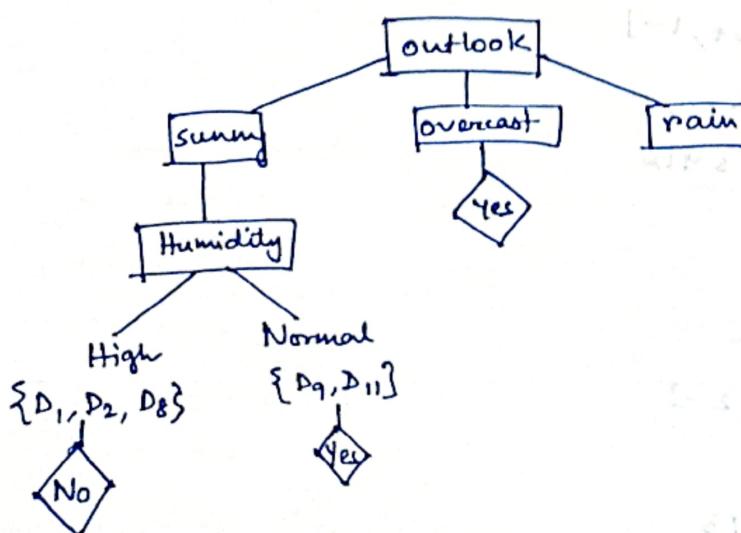
$$IG(\text{Sunny, humidity}) = 0.971 - \left(\frac{3}{5} \times 0 + \frac{2}{5} \times 0 \right) = 0.971 \rightarrow \max^m$$

wind

$$\rightarrow \text{Weak} \rightarrow [1+, 2-] = 0.918$$

$$\rightarrow \text{Strong} \rightarrow [1+, 1-] = 1.0$$

$$IG(\text{Sunny, wind}) = 0.971 - \left(\frac{3}{5} \times 0.918 + \frac{2}{5} \times 1.0 \right) = 0.0192$$



Partial Decision Tree

for rain

$$S(\text{Rain}) = [3+, 2-] = 0.971$$

temperature

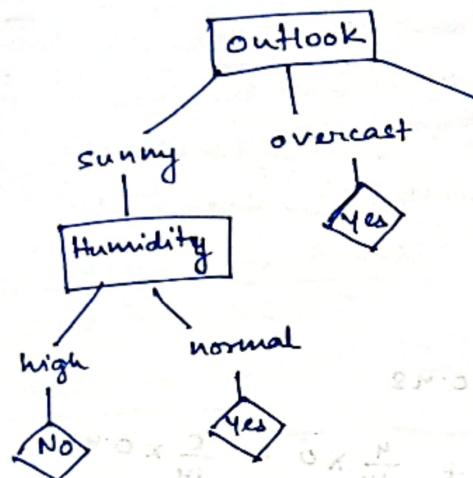
$$\begin{aligned} \rightarrow S_{\text{hot}} [0+, 0-] &= 0 \\ \rightarrow S_{\text{mild}} [2+, 1-] &= 0.9183 \\ \rightarrow S_{\text{cold}} [1+, 1-] &= 1.0 \end{aligned}$$

$$IG(\text{rain, temp}) = \frac{0.971}{0.0192} = 50.5$$

wind

$$\begin{aligned} \rightarrow S_{\text{weak}} \\ \rightarrow S_{\text{strong}} \end{aligned}$$

$$IG(\text{rain, wind}) = 0.971$$



Decision Tree

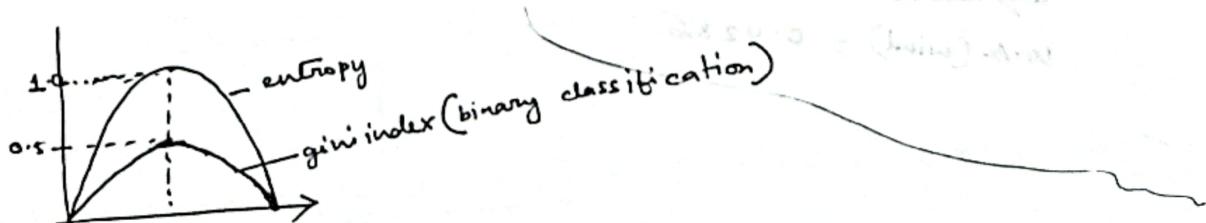
Gini Index/Gini Impurity/Gini coefficient

The another way of splitting the decision tree is via the Gini index. The entropy & IG method focusses on purity & impurity in a node. The Gini index measures the probability for a random instance being misclassified when chosen randomly. The lower the gini index, the better the likelihood of misclassification.

$$G_i = 1 - \sum_{i=1}^c p_i^2$$

where $p_i \rightarrow$ proportion of data points belonging to class i in the dataset.

- (i) For the context of binary classes, G_i has maximum value of 0.5 & min value of 0
- (ii) For multiple classes, G_i has value from 0 to 1.



Q. Decision Tree dataset

Gini index of ~~entire~~ dataset:

9 instances \rightarrow YES, 5 instances \rightarrow No

$$G_i(S) = 1 - \left[\left(\frac{9}{14} \right)^2 + \left(\frac{5}{14} \right)^2 \right]$$

$$\approx 0.45919$$

Computation of Gini index

Outlook

- \rightarrow sunny (5 instances) $[2+, 3-]$
- \rightarrow overcast (4 ") $[4+, 0-]$
- \rightarrow rain (5 ") $[3+, 2-]$

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

$$Gini(S_{\text{overcast}}) \rightarrow 1 - \left(\frac{4}{4} \right)^2 = 0$$

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

$$\therefore \text{weighted average (Outlook)} = \frac{5}{14} \times 0.48 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.48 = 0.34285 \text{ min.} \rightarrow \text{root}$$

Temp

- \rightarrow hot
- \rightarrow cool
- \rightarrow mild

$$\text{weighted avg (temp)} = 0.4405$$

Humidity

- \rightarrow high
- \rightarrow normal

$$\text{w.A. (humidity)} = 0.3673$$

Wind

- \rightarrow strong
- \rightarrow weak

$$\text{w.A. (wind)} = 0.4286$$

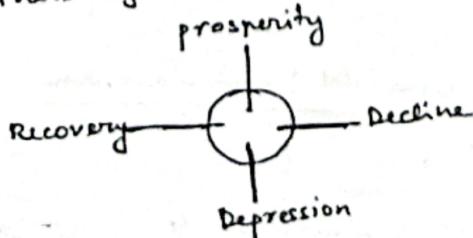
MODULE-3

TIME SERIES ANALYSIS

- A series of observations recorded in accordance with the time of occurrence is called series. The values of a variable are observed chronologically in terms of days, week, months, or year.
- The objects of time series analysis are to isolate and measure the effects of various components. Such analysis will help in understanding the past behaviour so that we may be able to predict the future tendencies.
- 4 components of time series
 - (i) Secular Trend / Trend (T)
 - (ii) Seasonal variations (S)
 - (iii) Cyclical Fluctuations (C)
 - (iv) Irregular / random movement (I)
- multiplicative model
$$y_i = S \times T \times C \times I$$
- additive model
$$y_i = S + T + C + I$$

VARIOUS COMPONENTS OF TIME SERIES

- (a) Secular Trend (T) \rightarrow It is the smooth, regular and long term movement exhibiting the growth/decline over a period of time.
Eg- population growth, death rate
- (b) Seasonal Variation (S) \rightarrow It represents a type of periodic movement where the period is not longer than 1 year. This up and down movement of timeseries recurring with remarkable regularity year after year, is attributable to the presence of seasonal variations.
- (c) Cyclical Fluctuations (C) \rightarrow It is another type of periodic movement where the period is more than a year. Such movements are regular & oscillatory in nature.
Eg- business cycle.



- (d) Irregular / random movement (I) \rightarrow This variations are caused by factors of an erratic nature. This variations don't follow a particular model, & so they are not predictable.

- increasing demand of automobiles — secular trend.
- passenger traffic during 24 hours a day — seasonal fluctuations
- recession — cyclic fluctuation
- a fire in a factory delaying production for a week — irregular movement

MEASUREMENT OF TREND

There are four methods of isolating secular trend in time series.

1. Free hand method.

2. Semi-average method

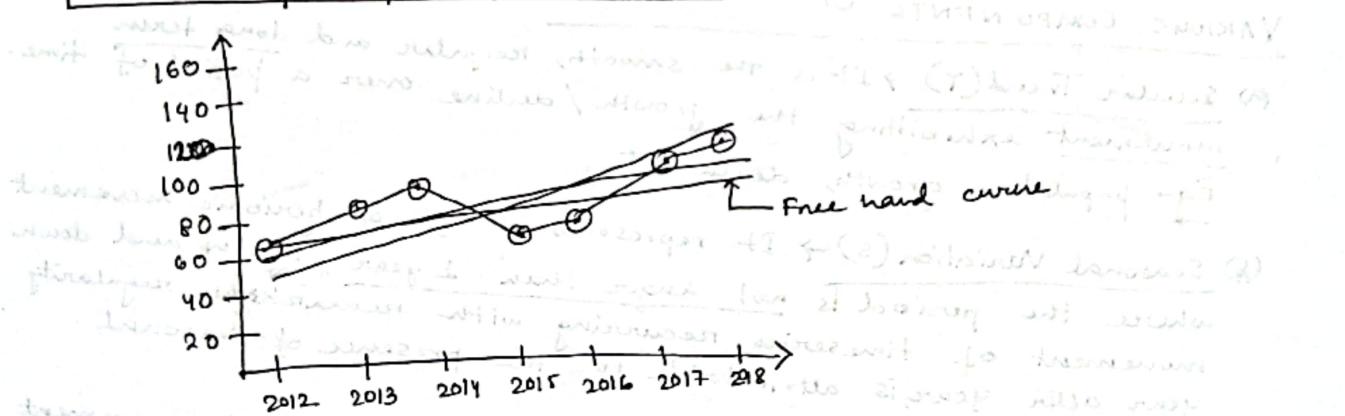
3. moving-average "

4. fitting mathematic curves.

1. Free hand method —

- The given data are plotted as points on a graph paper against time. The time-series data are shown along the vertical axis & time (t) along the horizontal axis.
- A smooth free-hand curve is drawn through the scatter of the plotted points. The dist. of this line is known as trend line which gives the trend value for each time period.

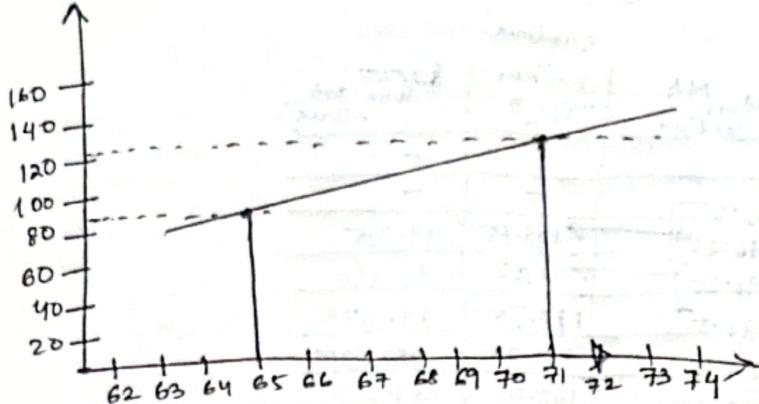
Time	2012	2013	2014	2015	2016	2017	2018
Values	64	82	97	71	78	112	115



2. Semi-Average Method —

The method consist in dividing the data into 2 parts and then finding an average of each part. These averages are plotted as points on a graph paper against the meet point of the time interval covered by each part. The st. line joining these 2 points gives the trend line.

year	time series	7 year semi total	semi-averages
$n=13$	1962 64		
	63		
	97		
	65	619	88.4
	66		$88.4 + 5.5 = 93.9$
odd — 1 year common	67 112		$93.9 + 5.5 = 99.4$
	68 115		
	69 131		
	70 88	850	121.4
	71 100		
	72 146		
	73 150		
		79120	



$$\begin{aligned}
 & \frac{121.4 - 88.4}{6} \\
 & = \frac{33.0}{6} \\
 & = 5.5 \text{ slope.}
 \end{aligned}$$

method of selected points.

It gets affected by outliers.

Ex - 2
Even dataset

Year	production	4 year semi-total	Semi avg
1971	40		
72	45	167	41.75
73	40		
74	42		
75	46	215	53.75
76	52		
77	56		
78	61		

3. method of moving average

- (i) A very convenient tool for ironing out fluctuations in time-series.
- (ii) It tends to reduce the amount of variations in the data sets, eliminating unwanted fluctuations (smoothing of time series).
- (iii) MA of order 'n' = $\frac{y_1 + y_2 + y_3 + \dots + y_n}{n}, \frac{y_2 + y_3 + \dots + y_{n+1}}{n}, \frac{y_3 + y_4 + \dots + y_{n+2}}{n}, \dots$

Note: (i) If the period of moving avg is odd, the trend values correspond to the given time.

(ii) If the period is even, a two-point moving avg of the moving averages so obtained has to be found out for centering them.

Ex - 1
calculate the 5-year (odd) moving avg

Year	value	5 year moving total	Syr MA
2000	105	-	-
01	115	-	-
02	107	490	98
03	90	480	96
04	80	450	90
05	95	425	85
06	85	395	79
07	75	-	-
08	60	-	-

Eg-②

Calculate the 4 year MA

Year	value	4-yr Moving total (not centered)	4 year MA (centered)	2-item moving total	4-yr MA (centered)
2000	105				
01	115	410	102.5		
02	100	385	96.25	198.75	99.375
03	90	365	91.25	187.5	93.750
04	80	350	87.5	178.75	89.375
05	95	335	83.75	171.25	85.625
06	85	315	78.75	162.5	81.250
07	75				
08	60				

merits

- The method is simple to apply & involves no difficult calculations.
- inf. • If the time series contains regular cyclical fluctuations, these fluctuations are automatically removed provided an appropriate period is chosen. Even when the fluctuations are not completely eliminated, MA process reduces the intensity.
- This method is flexible in the sense that if some more observations are added to the original series, the entire calculations need not to be changed.

Demerit

- Trend values for all the given time periods can't be obtained. Some trend values at the beginning & at the end of the series have to be left out. Their no. increases with inc. in the period of MA.
- The method can't be used for forecasting future trends
- MA are useful only when the trend is linear. If the graph is curve-linear, MA may deviate from the trend.

3(b) Weighted MA

Q. Find the trend of the following series using a 3 year weighted MA.

with weight 1, 2, 1

Year	values	3-yr moving total	3 year MA
1	2	—	—
2	4	$(2 \times 1) + (4 \times 2) + (5 \times 1) = 15$	3.75 ←
3	5	$(4 + 10 + 7) = 21$	5.25
4	7	27	6.75
5	8	33	8.25
6	10	41	10.25
7	13	—	—

$$15 / (1+2+1) = 15/4$$

4) Method of fitting mathematical curves

In this method, we use

[ONLINE CLASS NOTES IN LAPTOP]

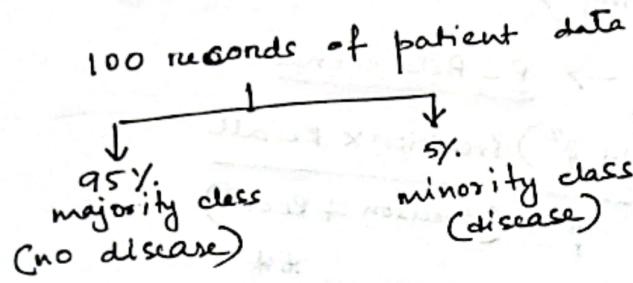
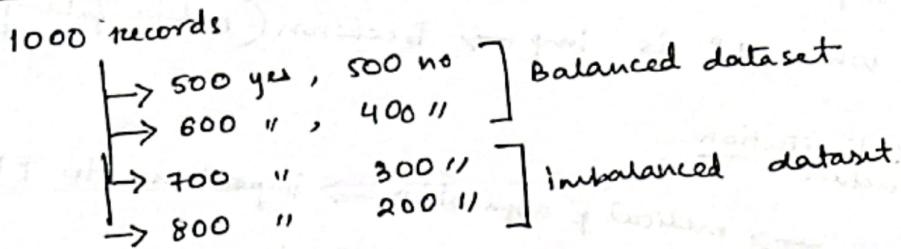
Method of least squares

Normal Eqns

UNIT - 4

CLASS IMBALANCE PROBLEM

Classification problem

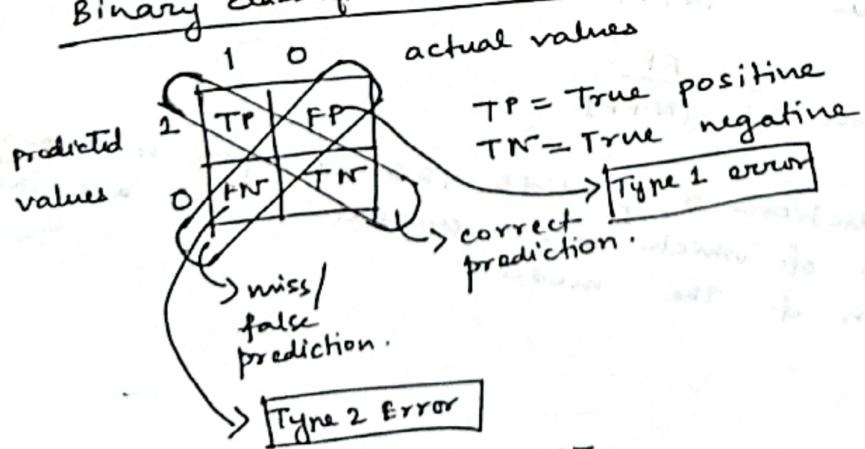


* classifier biased towards no disease (majority class)

CONFUSION MATRIX

In ML, to measure the performance of the classification model we use the confusion matrix. It is a matrix that summarises the performance of ML model on a set of test data. It is a means of displaying the no. of accurate & inaccurate instances based on the model's predictions.

Binary classification model



FP = False +ve
FN = " -ve

$$\text{accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

$$\text{Error rate} = 1 - \text{accuracy}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

or
TPR
(True Positive Rate)

or
Sensitivity

Mail Spam detection

		Actual	
		0	1
predicted	0	FP	importance to FP importance to correct unwanted mails.
	1	FN	when FP is imp \rightarrow Precision (reduce false positive)

Cancer detection

		Actual	
		0	1
predicted	0	FP	medical application \rightarrow importance to FN
	1	FN	when FN imp \rightarrow use Recall.

when both FN & FP is imp \rightarrow F- β score

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

if $\beta=1 \Rightarrow$ F-1 Score

$$= \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (\text{Harmonic mean})$$

more value
better model.

$\checkmark \text{TPR}$
(True positive rate) $= \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{\text{TP}}{\text{Actual positive}}$

$\checkmark \text{FNR}$
(False negative rate) $= \frac{\text{FN}}{\text{actual true}} = \frac{\text{FN}}{\text{TP} + \text{FN}}$

$\checkmark \text{TNR}$
(True negative rate) $= \frac{\text{TN}}{\text{actual -ve}} = \frac{\text{TN}}{\text{TN} + \text{FP}}$

$\checkmark \text{FPR}$
(False true rate) $= \frac{\text{FP}}{\text{actual -ve}} = \frac{\text{FP}}{\text{TN} + \text{FP}}$

Q. A model makes predictions & predicts 120 examples as belonging to minority class, 90 of which are correct & 30 are incorrect. What is the precision of the model.

Sol: precision $= \frac{90}{120} = 0.75$

Q.

		predicted	
		no	yes
Actual	no	50	10
	yes	5	100

Find accuracy, precision, recall.

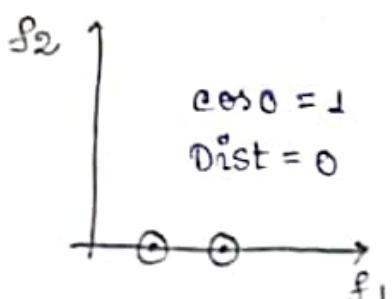
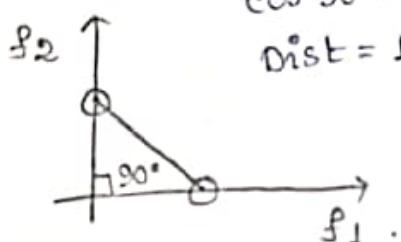
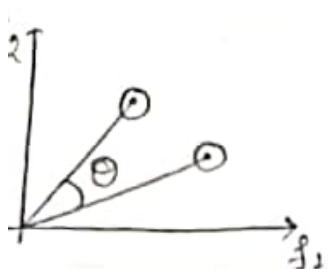
Ans:-

$$\begin{aligned}
 \text{Accuracy} &= \frac{TP + TN}{TP + TN + FP + FN} \\
 &= \frac{100 + 50}{50 + 10 + 5 + 100} = \frac{150}{165} = 0.909 \approx 0.91
 \end{aligned}$$

$$\begin{aligned}
 \text{Precision} &= \frac{TP}{TP + FP} & \text{Recall} &= \frac{100}{100 + 5} \\
 &= \frac{100}{100 + 10} & &= 0.9523 \approx 0.95
 \end{aligned}$$

Cosine similarity $\rightarrow \cos \theta$

Cosine distance $= 1 - \cos \theta$



$$\cos 90^\circ = 0 \quad (\text{similarity} = 0)$$

Dist = 1.

$$\cos(\alpha, \gamma) = \frac{\alpha \cdot \gamma}{\|\alpha\| \|\gamma\|}$$

dot product of vectors α, γ
length of the vector α, γ

cosine similarity is a metric used to measure the similarity of two vectors. Specifically it measures the similarity in direction or orientation of vectors ignoring differences in their magnitude. The similarity of two vectors is measured by the cosine of the angle b/w them.

Q. Find the cosine similarity b/w 2 vectors x & y .

$$x = \{3, 2, 0, 5\}$$

$$y = \{1, 0, 0, 0\}$$

$$\text{Ans: } x \cdot y = 3 \times 1 + 2 \times 0 + 0 \times 0 + 5 \times 0 = 3.$$

$$\|x\| = \sqrt{3^2 + 2^2 + 0^2 + 5^2} \\ = \sqrt{13 + 25} \\ = \sqrt{38} \\ = 6.16$$

$$\|y\| = \sqrt{1^2 + 0^2 + 0^2 + 0^2} \\ = 1.$$

$$\cos(x, y) = \frac{3}{6.16 \times 1} = 0.482 \approx 0.49.$$

$$\text{Disimilarity of } \cos(x, y) = 1 - 0.49 = 0.51.$$

Naive Bayes classifier.

Another name of mutual exhaustive set are called sample.

It is a supervised learning algorithm which is based on Bayes Theorem and used for solving classification problem

It is a probabilistic classifier which means it predicts on the basis of probability of an object.

The algorithm is comprised of two words \rightarrow Naive & Bayes

It is called Naive because it assumes the occurrence of certain feature is independent of occurrence of other features.

Bayes: - It follows the principle of Bayes theorem.

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)} \quad \begin{array}{l} \text{Joint probability} \\ \text{Marginal probability} \end{array}$$

Q. Determine whether the car is stolen or not for the ~~unseen~~ Sample X using Naive Bayes Method.

Sample x using Naive Bayes method.

An Unseen Sample

X = med, sporty, Domestic

Target class level: stolen.

□ Union prob :-

$$P(\text{yes}) = 5/10 = 0.5$$

$$P(\text{zero}) = 5/10 = 0.5$$

Likelihood for colon:-

$$P(\text{colon} | \text{yes}) \neq P(\text{colon} | \text{no})$$

Type (color) = Red & yellow.

$$P(\text{Red}/\text{yes}) = \frac{P(\text{Red} \cap \text{yes})}{P(\text{yes})} = \frac{3}{5}$$

$$P(\text{Red}/\text{No}) = \frac{P(\text{Red} \cap \text{No})}{P(\text{No})} = \frac{2}{5}$$

$$P(\text{Yellow} | \text{Yes}) = \frac{P(\text{Yellow} \cap \text{Yes})}{P(\text{Yes})} = \frac{2}{5}$$

$$P(\text{Yellow/No}) = \frac{3}{5}$$

Likelihood for Type :-

$P(\text{Type/No}) \neq P(\text{Type/Yes})$

$$P(\text{spont}/\text{no}) = 2/5 \quad P(\text{spont}/\text{yes}) = 4/5$$

$$P(\text{Surv}/\text{No}) = 3/5 \quad P(\text{Surv}/\text{Yes}) = 1/5$$

□ Likelihood for origin .

$P(\text{origin/No}) \neq P(\text{origin/Yes})$.

$$P(\text{domestic/No}) = \frac{3}{5} \quad P(\text{domestic/Yes}) = \frac{2}{5}$$

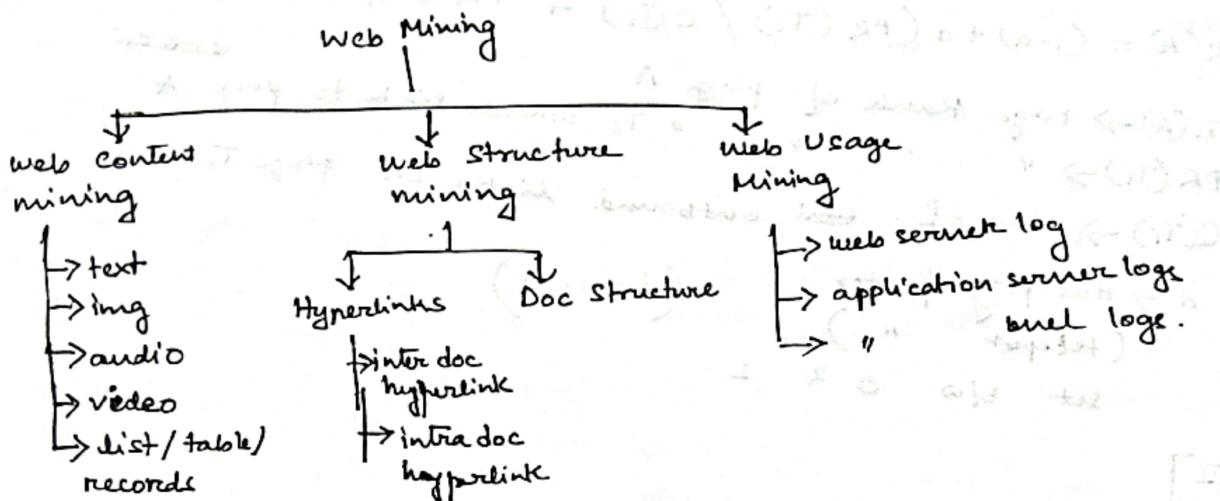
$$P(\text{imported/No}) = \frac{2}{5} \quad P(\text{imported/Yes}) = \frac{3}{5}$$

WEB MINING

{ Page rank algo
Hits algo

- (i) It is the art & science of discovering patterns & insights from world wide web (www) so as to improve it.
- (ii) It analyses data from the web and helps find insights that would optimize the web content & improve UX.

Challenges of mining www data. (Note → classroom)



Web content mining - (i) It is the process of extracting useful info from the contents of web doc. The content comprises of text, img, audio, video or ~~structured~~ records, etc.

(ii) The text & application contents on the pages could be analyzed for its usage by visit counts. It depends on the quality of content that attracts more users.

Web structure mining - (i) It is the process of discovering structured info (linked structures) from the web. The web walks through a system of hyperlinks using HTTP. Any page can create a hyperlink to any other page.

(ii) web-server data includes IP addresses, page reference, access time

Q. what is meant by authoritative web-page?
→ The web consist not only of pages but also of hyperlinks pointing from one page to another. These hyperlinks contain an enormous amount of latent human annotations that can help automatically infer the notion of authority.

There are 2 basic ~~stage~~ in strategic models for websites.
① Hubs → these are pages with a large no. of interesting links. They serve as a gathering point where people visits, to accm a variety of info

(ii) Authority - People would gravitate towards pages that provides the most complete & authoritative info on a particular subject. These websites have the most no. of inbound links from other websites.

Page Rank Algo

- (i) method of web structure mining.
 - (ii) It helps in measuring relative importance of a web page, within a set of similar entities. The websites with more no. of links or more links from higher quality websites will be ranked higher.
 - (iii) It is an algo used by Google Search to rank websites in their search engine.

$$PR(A) = (1-d) + d \left(\frac{PR(T_i)}{c(T_i)} + \frac{PR(T_n)}{c(T_n)} \right)$$

PR(A) \rightarrow Page Rank of page A
"T: which link to page A

PR(T_i) \rightarrow "page T_i "

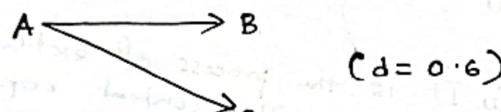
$C(T_i) \rightarrow$ no. of ~~out~~ outbound links on page i

$d \rightarrow$ damping factor $(d = 0.85)$

(tekeport ")

set b/w 0 & 1

Part II
969



initially page rank of all $\overset{C}{\underset{[]}{\text{pages}}} = 1$ (iteration 0)

$$\begin{aligned} \text{PR}(A) &= (1-0.6) + 0.6 [0 + 0] \\ &= 0.4 \end{aligned}$$

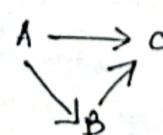
$$PR(B) = (1-0.6) + 0.6 \left[\frac{0.4}{2} \right] = 0.52$$

$$PR(C) = (1-0.6) + (0.6) \left[\frac{0.4}{2} \right] \\ = 0.52$$

$$\text{For } PR(A) = (1 - 0.4) + 0.6 \times 0 = 0.6$$

$$\text{iter}^{-2} \quad \text{PR}(B) = 0.52$$

PR (c) = 11



$$d = 0.85$$

$$PR(A) = PR(B) = PR(C) = 1$$

iter = 1

$$P(A) = (1 - 0.85) + 0.85 \times 0$$

$$= 0.15$$

$$PR(B) = (1 - 0.85) + 0.85 \left(\frac{0.15}{2} \right)$$

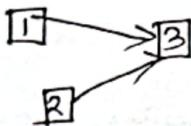
$$= 0.2137$$

HITS Algo

HITS - Hyperlink Induced Topic Search.
i) It is a linked analysis algo that rates web pages as being hubs / authorities.

ii) It uses hubs & authorities to define a recursive relationship b/w web pages.

Q. complete the hub & authority weights for the following graph.



S1 find the ~~adj~~ matrix of web pages. If there is a link from page i to page j, then matrix[i][j] = 1, else 0.

$$A = \begin{bmatrix} 1 & 0 & 1 \\ 2 & 0 & 0 \\ 3 & 0 & 0 \end{bmatrix}$$

S2 find the matrix transpose

$$A^t = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 0 & 0 \\ 2 & 0 & 0 \\ 3 & 1 & 1 \end{bmatrix}$$

S3 assume initial hub weight vector
compute the authority wt. vector

$$\begin{aligned} v &= A^t \cdot u \\ &= \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ 0 \\ 2 \end{bmatrix} \end{aligned}$$

S4 update hub weight

$$\begin{aligned} u &= A v \\ &= \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} 0 \\ 0 \\ 2 \end{bmatrix} \\ &= \begin{bmatrix} 2 \\ 2 \\ 0 \end{bmatrix} \end{aligned}$$

Node 1 \rightarrow hub
2 \rightarrow hub
3 \rightarrow authority.

Regression Analysis

- (i) Statistical method to model the relationship b/w a dependent (target) & independent (predictor) variables with one/more independent variables.

$$y = f(x)$$

↓
dependent independent

- (ii) It predicts continuous/real values. (eg. Temperature, age, salary, etc.)
- (iii) It is supervised learning technique which helps in finding the co-relation b/w variables & enables us to predict the continuous o/p variable.
(Correlation coeff = -1 to +1)
- (iv) Mainly used for prediction, forecasting, time-series modelling & determining the causal effect relationship b/w variables.
- (v) In regression, we plot a graph, which best fits the given data point, & using this plot the ML model makes prediction
- (vi) The main factor in regression which we want to predict is called the dependent variable (target variable)
- (vii) The factors which affect the dependent variables / which are used to predict the values of dependent variables are called independent variables (predictor).

Error in Simple Regression

The regression eq. model in ML uses the slope intercept format. It identifies the values of intercept & slope by relating the values of x & y . However identifying the exact match of values for intercept & slope is not always possible. There will be some error (ϵ) associated with it. ϵ is called marginal/residual error. (Residual is the dist b/w predicted point on ~~act~~ regression line & actual point)

(ix) A st. line is drawn as close as possible over the points on the scatter plot. Ordinary Least Square (OLS) is the technique used to estimate a line that will minimise the error (ϵ), which the difference b/w the predicted & actual ~~values~~ of y . This means summing the error of each prediction, (sum of the sq. of the errors)

$$SSE = \sum_i \epsilon_i^2$$

$$\text{if } y = (a + bx) + \epsilon$$

It is observed that SSE is least when 'b' takes

$$b = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

The corresponding value of 'a' is calculated as

$$a = \bar{y} - b \bar{x}$$

OLS algorithm

- S1: calculate the mean of x & y
- S2: calculate the errors of x & y
- S3: obtain the product
- S4: obtain the sum of the products
- S5: sq. the diff. of x

S6: Get the sum of the sq. difference

S7: Divide o/p of S4 by o/p of S6 to calculate slope 'b'

S8: Divide o/p of S4 by o/p of S6 to calculate slope 'b'

- Q. A random sample of 15 students in a class was selected & internal & external exam data is given below.
- Construct the regression line.
 - What is your estimate of the marks a student can obtain in the ext. exam if the student obtain 25 in internal.

X	Y	$X - \bar{X}$	$Y - \bar{Y}$	$(X - \bar{X})^2$	$(X - \bar{X})(Y - \bar{Y})$
15	49				
23	63				
18	58				
23	60				
24	58				
22	61				
22	60				
19	63				
19	60				
16	52				
24	62				
11	30				
24	59				
16	49				
23	68				
Σ	299	852		226.9335	429.8

$$\bar{x} = 19.933$$

$$\bar{y} = 56.8$$

$$y = a + bx$$

$$b = \frac{429.8}{226.9335} = 1.894$$

$$\therefore y = 19.047 + 1.894x$$

$$(ii) 66.397$$

$$a = \bar{y} - b\bar{x}$$

$$= 56.8 - 1.894 \times 19.933$$

$$= 19.047$$