

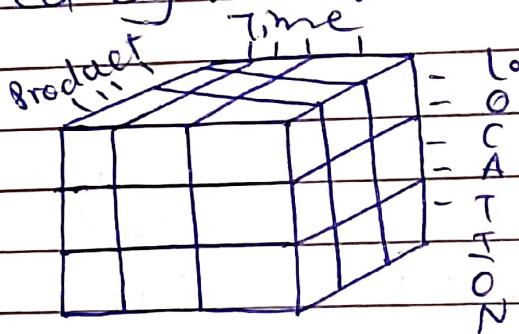
- Reporting Authoring -

* Multidimensional Data Model - (MDM)

MDM is a way of organizing and representing data that focuses on multiple aspects or dimensions of information particularly useful for analyzing large datasets involving various factors and allows for faster retrieval.

> OLAP and DW uses MDM databases. It represents data in the form of data cubes. Data cubes allow to ~~model~~ model and view data from many dimensions & perspectives.

> It is defined by dimensions and facts and is represented by fact table



> Once Data cube is ready it can be queried using OLAP such as Slicing, dicing, drillup, drilldown, rollup, Pivoting

* Relational Data Model (RDM) -

RDM Structured way of organizing data into tables (relations) consisting of rows & columns. Data is stored and accessed and managed using query language SQL

Key components -

> Tables - table has name & consists of rows (records) and column (attributes).

- 2) Rows - Individual records in a table, representing a single instance of entity.
- 3) Columns (Attribute) - Properties of entity
- 4) Primary key - Unique identifier for each row within table.
- 5) Foreign key - These are columns that references primary key of another table.
Ex - Employee table.

ID	Name	Email	Dept ID	Position
1	XYZ	--@gmai	101	Manager
2	ABC	--@gmai	102	Analyst
3	SYT	--@mai	103	Developer.

Department table -

Dept ID	Dept Name	Location
101	Sales	NYC
102	Marketing	IND
103	IT	BAN

Employees table has foreign key (Dept ID) that establishes a relationship with Department table.

* Types of Reports -

- 1) List reports - Present detailed info in a tabular format often with rows and columns. Each row represents a data point each column represents an attribute.

Use case

- > Displaying customer lists, product catalog

2) Crosstab Reports - Analyze data by summarizing relationships between 2 or more categorical variables in matrix format. Rows represents variables with categories on the rows and columns and intersection of these categories displays the corresponding value (frequency, count sum, etc.).

Use case

- 1) Identify trend & patterns b/w variables.
- 2) Analyze the relationship b/w variables.

3) Statistical Reports - Summarizes & analyzes quantitative data using statistical measures like mean, median, mode, variance etc. may include tables and charts to present measures.

Use case

- 1) Providing high quality overviews of central tendency & dispersion.
- 2) comparing different dataset using statistical metrics.

4) Chart Reports - visually represent data patterns and trends using charts like bar charts, line charts, pie charts or Scatter plots. Provides visual pattern & relationships.

Use cases

- 1) Highlighting trends, relations between datapoints.
- 2) Visual representation.

3) Map Reports - Display data geographically using maps. Can be helpful for analyzing trends, distributions or patterns with location.

use case -

1) Identifying geographic areas or specific patterns.

2) Understanding the spatial context of your data.

4) Financial Reports - Present financial info about company or organization. They typically include income statements, balance sheet, cash flow etc.

* Data grouping & sorting -

> Grouping - Involves organizing data points into categories based on shared characteristics.

Grouping can make large datasets more manageable by focusing on specific or subset of data.

Grouping can be used for summarizing data by categories, analyze data pattern & trends in group.

> Sorting - Arranging data in specific order, either ascending or descending, based on one or more columns. Sorting is crucial for making it easier to read & analyze.

Sorting used for organizing data in logical sequence, To improve readability, improved data analysis.

* **Filtering Reports** - It is the process of refining data to display only the relevant info that meets specific criteria. It helps users focus on subsets of data for analysis, making reports more insightful.

> Benefits

- 1) Focus on specific Trends
- 2) Identify outliers
- 3) Simplify complex data.

> Types of Filters

- 1) Value filters - based on specific values in one or more column
- 2) Range filters - based on range of values.
- 3) Text filters - filtering based on text patterns.
- 4) Date filters

* **Drill Down** - (zooming in) -

Drill down allows user to navigate from higher level to lower level of details within the data. It helps users to break down aggregated data into finer details, uncovering the underlying component.

Example - Imagine sales report that shows Sales by region (eg North, East, West..). Drilling down on the North region might reveal sales figures for each state within North region.

* **Drill up** (zoom out) - opposite to drill down.

It allows user to navigate from a lower level of details in data to a higher level of aggregation. It collapses detailed view for a more concise view.

Ex - Continuing with sales report drilled down at North region user could drill up to go back to original view showing total sales by region.

* Drill through - Allows user to navigate from a summarized view in a report to a completely different but related data.

This action is typically initiated by clicking on a specific element within the report it takes user to a new view

Ex- Imagine a report showing total customer purchases by product category. Drilling through on a specific product category (eg. laptops) might take you to separate report that details individual laptop models sold, their prices, etc.

Drillthrough

Data preparation

Page No.	
Date	

* Incomplete data -

This term refers to datasets that have missing values or lack certain necessary attributes. This can happen due to various reasons such as data entry errors, loss of data.

> Implication -

1) Bias in Analysis - Missing data can lead to biased results if not handled properly.

2) Reduced statistical power - Incomplete datasets can weaken statistical power of an analysis, making it harder to detect true patterns.

3) Inaccurate Model - ML models trained on incomplete data may not generalize well and can produce unreliable predictions.

> Strategies to handle incomplete data -

1) Data manipulation - Filling in missing values

1) Mean, median, mode imputation - Replacing missing values using the mean, median, mode of column.

2) KNN algorithm (KNN) - Using the nearest neighbors to estimate the missing values.

2) Deletion method - If missing data is significant or cannot be reliably imputed, removing affected records from records might be necessary.

* Data Affected by Noise -

Noisy data refers to datasets that contain errors, outliers or irrelevant data points that do not represent true values or patterns.

ML model might learn the noise instead of the underlying pattern, reducing their ^{noise} instead of the underlying, reducing performance on new data. (generalization).

> Strategies -

- > Data Cleaning - Involves identifying and removing noisy data or replacing them using mathematical techniques like Z-scores, standardization, calculating Inter Quartile range
- 2) Data filtering - Filtering out records with extreme values of those exceeding a certain threshold to reduce influence of outliers.
- 3) Clustering Algorithms - Using clustering algorithms like DBSCAN, K-Means can be used to find the main groups in the data and treat points far from any cluster center as noise.
- 4) PCA - Reduce dimensionality by transforming data into set of principal components that capture the most significant variance.
- 5) Regression techniques - Use regression techniques that are less sensitive to outliers such as RANSAC or Huber regression.

* Data Transformation - Transformer -

- > Standardization - Also known as Z-score normalization, involves transforming the data so that it has mean of 0 and standard deviation of 1

$$T_z = \frac{z - \mu}{\sigma}$$

Used in algos like that depends on distance metrics (eg KNN, SVM). In optimization algos like gradient descent for faster convergence

By putting features on similar scale, z score normalization ensures all features contribute equally during training.

→ z score normalization can be sensitive to outliers significantly affecting mean & standard deviation.

~~2) Min Max Scaling~~

~~→ Feature extraction~~

2) Min Max scaling - Also known as normalization, data transformer technique used that adjusts the range of the features to scale values to specific range, typically $[0, 1]$ or $[-1, 1]$.

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

→ Feature Extraction - Key step in data preprocessing pipeline for ML & data analysis focuses on identifying and extracting the most relevant and informative characteristics from raw data. Purpose -

- 1) Dimensionality reduction - reduces no. of features while retaining important info.
- 2) Improved Model performance
- 3) Noise reduction.

Common techniques

1) Principal Component Analysis (PCA) -

PCA is dimensionality reduction technique commonly used in data analysis. PCA aims to simplify data by identifying the underlying structure and reducing no. of features while preserving imp. data.

> PCA Process

1) Standardize data. - Z score normalization.

2) Covariance matrix - Compute covariance matrix, which captures relationships b/w features.

3) Eigen decomposition - Perform eigen decomposition on covariance matrix to find eigenvalues & eigenvectors. Eigenvector represents direction of highest variance & eigenvalue represents magnitude of variance.

4) Selecting Principal Component - Selecting top K eigenvectors corresponding to the rank by arranging eigenvalues in descending order.

5) Dimensionality reduction - by selecting top K principal components dimensionality reduction is achieved.

> Benefits

1) Reduced complexity

2) Improved ML performance.

3) Feature extraction.

➤ Data Reduction - Important process in data preprocessing that helps to manage large datasets by reducing their size while retaining original essential info

➤ Techniques

➤ Sampling - Involves selecting a subset of data points from the larger dataset. This subset should represent the overall dataset adequately, allowing for efficient analysis.

Types of Sampling -

i) Random Sampling - Selecting data points randomly from the entire dataset. Each data point has an equal probability of being selected.

These can be chosen in various ways by selecting each data point randomly or dividing dataset into distinct subgroups based on certain characteristic.

ii) Stratified Sampling - Data is split into partitions and samples are drawn from each partition.

iii) Systematic Sampling - Data points are selected at a fixed interval from an ordered list. This method is efficient but requires data to be well structured.

iv) Clustering Sample -

v) Cluster Sampling - Population divided into groups (clusters) based on shared characteristics. Then random samples are chosen from each cluster.

2) Feature Selection -

Involves Selecting the most relevant features (variables) from the dataset while discarding less important or redundant features.

Helps in improving model performance, reducing overfitting & decreasing computational cost.

3) Types -

i) Filter method - These techniques evaluate individual features based on statistical measures like correlation with target variable. Features with highest scores are selected. These methods are generally computationally efficient but might not capture complex relations. Ex - chi-square test, Correlation analysis, F-value, F-value Calculation For ANOVA.

ii) Wrapper method - Evaluate feature subsets by training and testing a machine learning model on different combinations of features. Aims to find the subset that produces the best model performance.

Techniques

i) Forward - starts with an empty set of features and adds features one by one.

ii) Backward Elimination - Starts with all features and removes them one by one evaluating model performance at each step.

3) Embedded methods - This method performs feature selection as part of model training. Some algorithms like LASSO regression or tree based methods performs feature selection during training by assigning weights or importance scores to features.

4) Data discretization - Also known as binning, it is the process of grouping continuous values of variables into contiguous intervals. It is a data reduction mechanism because it diminishes data from large domain of numeric values to a subset of categorical values.

> Steps -

- 1) Sorting continuous values of the features
- 2) Evaluating a cut point for splitting intervals for merging.
- 3) Splitting or merging intervals of continuous values

> Typical Methods of Data Discretization -

1) Binning - Techniques in Binning

2) Equal width Binning - Simplest method, where the entire range of the continuous features is divided into specific no. of bins of equal length.

Ex. if the range of features is 0 to 100 and we choose 5 bins would cover an interval of 20 units
[0-20] [20-40] [40-60] [60-80] [80-100]

ii) Equal-Frequency binning - Data is divided into bins containing approximately the same number of data points. However, the bin widths might not be equal, especially if the data is not evenly distributed.

3) Histogram Analysis - It is an unsupervised discretization technique because it does not use class info like binning. There are various partition rules used to define histograms. In equal width histogram, values are partitioned in equal size bins. Histograms are effective for data with multiple attributes to capture dependencies between them.

4) Cluster Analysis - Popular data discretization method. A clustering algo can be applied to discretize a numeric, Partitioning the values into clusters or groups based on Similarity Scores. Partitioning data into clusters. There are many choices of clustering definitions & clustering algos like K-means or K-medoids.

4) Correlation analysis - Supervised descretization, method. Also known as chi merge algorithm. It is performed recursively by finding best neighbouring ~~intervals~~ intervals that have similar distribution or classes and merge them.

* Data Exploration

> Univariate analysis - It is simplest form of data analysis that involves examining each variable in a dataset independently. Helps to understand the basic structure and characteristic of individual variable.

> Graphical analysis of categorical Attributes - Purpose to visualize the distribution & frequency of categorical variable

Techniques

> Bar chart - Display freq or proportion of categories. Each Bar shows category & height represent freq.

> Pie chart - Shows relative proportions of different categories. Each slice represents a category's proportion to the whole

> Line graph - Used to track changes over period of time

> Graphical analysis of Numerical attributes -

To understand distribution, central tendency & spread of numerical value.

Techniques

> Histogram - Displays freq distribution of numerical variable

> Box plot - Summarizes data using minimum, 1st quartile, median, 3rd quartile & maximum value

> Density plot - Shows data distribuⁿ in a smooth continuous curve, useful for identifying data's shape.

> Measures of Central tendency - To identify the center point of a numerical dataset.

Measures -

> Mean = $\frac{\sum n}{n} \frac{\sum x}{n}$

- 2) Median - middle value.
- 3) Mode - Most frequent value.

> Measure of dispersion of numerical Attributes -
To understand spread or variability in dataset

Measures -

- 1) Range - Range = Max - Min
- 2) Variance - Avg of Squared difference from mean

$$\text{Variance } (\sigma^2) = \frac{\sum (x - \bar{x})^2}{n}$$

- 3) Standard deviation - Square root of variance indicating data deviaⁿ.

$$\text{Standard deviation } (\sigma) = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

- 4) Interquartile Range (IQR) - Difference b/w 3rd quartile & 1st quartile.

$$IQR = Q_3 - Q_1$$

> Identification of Outliers -

To detect datapoints that significantly differ from rest of dataset.

Technique

- 1) Box-Plot - Outliers as points beyond whiskers
- 2) Z-Score - $Z = \frac{x - \bar{x}}{\sigma}$ typically

$|Z| > 3$ are considered outliers.

- 3) IQR method.

2) Bivariate - Involves examining the relationship between two variables to understand how one variable influences or is associated with the other.

> Graphical Analysis -

To visually inspect relation b/w 2 variables.
Technique.

1) Scatter plot - used 2 numerical variables, showing data points plotted on 2D graph. The pattern indicates relations (positive, negative or no reln).

2) Line plot - useful for some time series data showing trends over time.

3) Box plot comparison - comparing numerical variable across different categories of variables.

4) Heatmap - Represents relations b/w 2 variables using colors to show co-relation.

> Measures of Correlation for Numerical Attribute
To quantify the strength & direction of relationship b/w 2 variables.

Technique

1) Pearson Correlation Coefficient (r) - measures linear relationship b/w two continuous variable
value range: -1 to +1

Positive value - Positive linear relation

Negative value - Negative linear relation

2) Spearman Rank Correlat - Measures monotonic relationship b/w two variables based on their rank

useful when data is not normally distributed.

3) Contingency table - To summarize relationship between 2 variables. A matrix that displays the freq. distribution of variables. Each cell represents the count of occurrences.

Chi-Square test used to determine significant association b/w 2 variables.

3) Multivariate - Involves examining more than 2 variables to understand their relationships.

This type reveal complex patterns in data.

→ Graphical Analysis -

1) Scatter plot matrix.

2) Heatmap

3) 3D scatter plot

4) Parallel coordinate plot - Line represent

each data point and each axis corresponds to different variable.

2) Measure of correlation.

1) Correlation matrix - A table of correlation b/w multiple variables.

2) Pearson correlation

3) Partial correlation - extends Pearson's correlation by far it helps to isolate the correlation b/w 2 variables while controlling effect of another variable.

Univariate	Bivariate	Multivariate
→ single variable at a time	Relationship between 2 variables.	between multiple variables.
→ understand distribution (center, Mean, Median, mode)	Identify pattern correlations	Explore complex relationships, identify clusters.
3) freq tables. Measures of central tendency (mean, median, Mode) variance, Standard deviation	Scatter plot, correlation coefficients (Pearson Correlation, Spearman's rank)	Scatter plot matrix pair plots, 3D Scatter plot, Partial Correlation.
4) To understand basic characteristic of single variable	Measure & describe strength and direction of relationship	To understand interactions and combined effects of multiple variables.
3) Graphical techniques Histogram Bar charts Box plot	Scatter plot Box plot Heatmap	Scatter plot matrix Pair plots, 3D Scatter plots etc.

Randomly.

~~Classification - It is a supervised learning task that involves predicting the category or class of a given data point. Goal is~~

~~A Classification - It is a fundamental task in ML that involves assigning data points to specific categories such as classifying reviews to positive or negative.~~

~~A Classification problems - These problems involve dependent & independent variable dependent variables used for training model to predict dependant variable.~~

eg - Spam ham classification emails.

Types of classification problems -

1) Binary classification - Predicting one of two possible categories.

2) Multiclass - Predicting one of many possible categories.

eg - handwritten digit recognition.

~~A Evaluation of classification model -~~

1) Accuracy - $\frac{\text{No. of correct prediction}}{\text{Total no. of prediction}}$

2) Precision - Ratio of correctly predicted positive observation to the total predicted positives.

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

3) Recall = $\frac{TP}{TP + FN}$

4) F₁ Score = $2 \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$

5) Confusion matrix - Table describing performance of classification model

0	TP	FP
1	FN	TN
0)		

c) AUC-ROC (curve - graphical representation of TP rate versus FP rate at various threshold settings. AUC represents Area Under the Curve.

~~A~~ Bayesian methods -

Y Bayes Theorem - Statistical and probability theory is a mathematical formula used to determine the conditional probability of events. It describes the probability of an event based on prior knowledge of the conditions that might be relevant to the event.

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

$P(A|B)$ (Posterior) - Probab. $P(A)$ after observing B
 $P(B|A)$ (Likelihood) - $P(B)$ given A is true
 $P(A)$ - $P(A)$ before observing B
 $P(B)$ - total probability of observing B.

It is based on 3 components.

- 1) Background knowledge - model's parameter
- 2) Data - observed evidences
- 3) posterior inference - combination of first 2 component

> Naive Bayes classifier - It is a classifier based on Bayes' Theorem with an independent assumption that all features are conditionally independent. Commonly used in text classification tasks like spam detection, sentiment analysis.

> Bayesian inference - used in fields like medical diagnosis, where the likelihood of a disease is updated based on test results.

ex - PFD

$$P(\text{Disease}|\text{PositiveTest}) = \frac{P(\text{PositiveTest}|\text{Disease}) \cdot P(\text{Disease})}{P(\text{PositiveTest})}$$

> Bayesian networks - Graphical model that represents probabilistic relationships among a set of variables

> Markov Chain Monte Carlo - Probability distribution based on constructing a Markov chain.

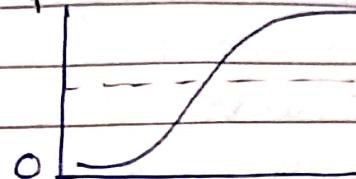
? Benefits of Bayesian methods -

- 1) Incorporating prior knowledge
- 2) Handling uncertainty
- 3) Flexibility - They can apply to problems with limited data or complex models.

(LR)

- > Logistic Regression - Fundamental ML task.
For classification tasks in ML, logistic regression deals with categorical dependent variables.
- > Logistic regression employs a Sigmoid function to map the linear combination of its features to probability value betⁿ 0 & 1.

$$\text{LR} = \frac{1}{1 + e^{-z}}$$



- > LR is trained using an iterative algo that aims to minimize Cost function.
- > By adjusting the weights and bias term through optimization process.
- > Application -
Spam filtering
Fraud detection
Customer churn.

Clustering - Grouping data points together based on their inherent similarities without predefined labels. It's like sorting a pile of fruits based on their characteristics (color, shape, etc)

- > Clustering methods -
- > Partitioning method - dividing data points into a fixed no. of clusters (K) in one go.
iteratively refine the cluster.
method

Partitioning models -

- 1) Kmeans clustering - widely used partition method. randomly assigns data points to k initial clusters and then iteratively perform steps -
 - > calculate centroid - mean value of each feature is calculated for each cluster creating centroid.
 - > Reassign Data points - Each data point is assigned to the cluster with nearest centroid.
 - > Update centroid & reassign.

2) K medoid clustering - similar to k mean but uses medoid instead of means. Medoid is centrally located data in cluster.

2) Hierarchical method - creates tree like structure representing relation between data points.

Models -

- 1) Agglomerative - starts by considering each data point as a separate cluster. Then iteratively merges the 2 closest clusters based on distance measures until single cluster remains.
- 2) Divisive - Approach starts with all data points in a single cluster and then iteratively splits cluster that maximizes a separability criterion until stopping criterion is reached.
- 3) Density - Based clustering method - type of unsupervised ML technique, density based method focus on areas with high concentration of data points

Models -

- 1) DBSCAN - Identifies high density areas by focusing on core points - data points surrounded by minimum number of neighbors within a specified radius.
Clusters are formed by connecting core points that are density reachable. ie reaching from one core to other core point through a chain of neighbors with radius.
- 2) OPTICS - Similar to DBSCAN, OPTICS focuses on density but takes different approach. It orders data based on reachability distance. No need to predefine the radius or minimum neighbors -

~~Evaluation of clustering models.~~

~~Clustering algos are~~

> Internal Evaluation Metrics.

1) Inertia (Within cluster sum of squares) -

It calculates sum squared distances between each data point and centroid of cluster.

Lower inertia implies that clusters are compact and well separated.

2) Silhouette Coeff - Evaluates how well data points are assigned to their clusters. It

considers both the distance to the assigned cluster's centroid and the distance to nearest neighbouring cluster's centroid.

- 3) Davies - Bouldin Index - Evaluate the average similarity ratio of each cluster with the cluster that is most similar to it. Lower values indicate better clusters.
- 4) Random Index - Measures the similarity betn clustering results and the ground truth labels. considers both true positives & true negatives.
- 5) Normalized Mutual Info (NMI) - Measures the amount of info shared between the clustering results and ground truth ~~bt~~ labels.
- 6) Calinski - Harabasz Index - Compares average distance between clusters to the variance within clusters. Higher value suggests better separation.
- 7) Davies Bouldin Index - It measures the ratio of the within-cluster scatter to the between cluster separaⁿ; lower value better clustering

* Association Rule - It is a technique in data mining that helps discover interesting relations between items in a large dataset. often applied in market basket analysis, where goal is to find associaⁿ between items purchased together.

Structure -

- 1) Antecedent (X) - set of item that appear in rule's premise (before the arrow (\Rightarrow)))
- 2) Consequent (Y) - set of items that appear in rule's conclusion (After the arrow (\Rightarrow)).
- 3) Support - Frequency of occurrence of both $X \& Y$
- 4) Confidence - Probability of Y appearing in transacⁿ given X is present

3) Lift - ratio of observed support to expected support if $X \& Y$ were independent.

> Apriori Algo - It is one of the algorithm used for transaction data in Association rule learning. Allows us to mine the frequent items in order to generate association rule b/w them.

- 1) Define min Support & min Threshold confidence
- 2) Find frequent item sets
- 3) Generate Association rules

Ex - TID	Items
T ₁	11, 12, 15
T ₂	12, 14
T ₃	12, 13
T ₄	11, 12, 14
T ₅	11, 13
T ₆	12, 13
T ₇	14, 13
T ₈	11, 12, 13, 15
T ₉	11, 12, 13

min support count is 2 ~~min count~~

min confidence is 60% Frequency Table -

TID	11	-	6
	12	-	7
	13	-	6
	14	-	2
	15	-	2

minimum support count is 2 so every item is kept in table no one is removed.

Frequency for item size 2

(11, 12)	- 4	(11, 12)
(11, 13)	- 4	(11, 13)
(11, 14)	- 1	(11, 15)
(11, 15)	- 2	(12, 13)
(12, 13)	- 4 min	(12, 14)
(12, 14)	- 2	$\xrightarrow{\text{Support}}$ (12, 15)
(12, 15)	- 2	
(13, 14)	- 0	
(13, 15)	- 1	
(14, 15)	- 0	

Frequency for item size 3

(11, 12, 13)	- 2 min.	(11, 12, 13)
(11, 12, 14)	- 1	$\xrightarrow{\text{Support}}$ (11, 12, 15)
(11, 12, 15)	- 2	

Association rule for (11, 12, 13)

min support 2 min confidence - 60%.

Rule	Support	Confidence
(11 ∩ 12) → 13	2	2/4 = 50%
(11 ∩ 13) → 12	2	2/4 = 50%
(12 ∩ 13) → 11	2	2/4 = 50%
13 → (11 ∩ 12)	2	2/6 ≈ 33.3%
12 → (11 ∩ 13)	2	2/7 ≈ 28.57%
11 → (12 ∩ 13)	2	2/8 ≈ 33.3%

Confidence = Support(A ∪ B) / Support(A).

$$\text{eg } \underbrace{(11 \cap 12)}_A \rightarrow \underbrace{13}_B = \frac{\text{Sup}(11 \cap 12) \cup 13}{\text{Sup}(11 \cap 12)} = \frac{2}{4} \leftarrow \begin{array}{l} \text{written support} \\ \text{of triplet formed.} \end{array}$$

\leftarrow unite support of 2(11, 12) together from freq table.

$$= \frac{1}{2} \approx 50\%.$$

∴ (11, 12, 13) does not hold any assoc' role.

For $(11, 12, 15)$

Rule	Support	Confidence.
$(11 \wedge 12) \rightarrow 15$	2	$2/4 = 50\%$. X
$(11 \wedge 15) \rightarrow 12$	2	$2/2 = 100\%$.
$(12 \wedge 15) \rightarrow 11$	2	$2/2 = 100\%$.
$11 \rightarrow (12 \wedge 15)$	2	$2/6 \approx 33.3\%$. X
$12 \rightarrow (11 \wedge 15)$	2	$2/7 \approx 28\%$. X
$15 \rightarrow (11 \wedge 12)$	2	$2/2 = 100\%$.

Consider those rules having confidence more than min confidence level ie 60%.
ie rules $(11 \wedge 15) \rightarrow 12$, $(12 \wedge 15) \rightarrow 11$ & $15 \rightarrow (11 \wedge 12)$

So the Association rule for $(11, 12, 15)$ is

$(11 \wedge 15) \rightarrow 12$, $(12 \wedge 15) \rightarrow 11$ & $15 \rightarrow (11 \wedge 12)$

Arounday

* Role of Analytical tools in BI -

Analytical tools are the backbone of BI. They act as powerful instrument for collecting, transforming, analyzing & visualizing data. to extract valuable insight for better decision making.

> BI tools can connect to various data sources together info

> Data might need cleaning, filtering & transformation to ensure consistency and usability for analysis.

> BI tools provide descriptive statistics and visualization to summarize data. By drilling down into data, analytical tools help identify the cause of observed outcome.

> Advanced Analytical tools help us to predict on new data using ml algos. This includes predicting Sales forecast, customer behaviour etc.

> Analytical tools allows creation of interactive dashboards. Tools enable the generation of detailed report that can be customized.

> After all analysis these tools helps in deeper further business decision making.

* Analytical tools -

> WEKA (Waikato Environment for Knowledge Analysis) -

Designed for data mining tasks provides collection of visualizing tools & algos.

> It offers wide range of ml algos, including classification, regression, clustering & association rule mining.

> Weka provides visualization tools to help understand the data and the data learning algo.

> It has user friendly GUI

- 2) KNIME - free open source platform for data mining & ml -
- > KNIME easily to use graph GUI to create workflows for data analysis, offers drag & drop workflow for building data analysis pipelines
 - > Set of powerful extensions and integrations make KNIME a versatile & scalable.
 - > Tools like R, Python, and WEKA can be integrated with KNIME.
 - > KNIME includes tools for data preprocessing, transformation, modeling & visualization.

- 3) RapidMiner - Data Science platform that provides an integrated environment for data preparation, ml, dl, text mining.
- > Offers a free open source edition and a paid commercial version.
 - > Similar to KNIME it uses drag and drop approach for building workflows.
 - > Includes data preprocessing, ml, text mining and data visualization.

- 4) R - freely available with vast collection of packages for various statistical & graphical techniques.
- > CRAN (Comprehensive R Archive Network) hosts thousands of packages that extends R's capabilities.
 - > R excels at data visualization, offering packages like ggplot2.

* Business Analytics - Science of extracting meaningful info from data to improve business decision making. BA focuses on transforming raw data into from various sources into actionable insights that can guide business strategies.

> Key Stages

- 1) Data acquisition - Data collect from various sources
- 2) Data analysis & Exploratn - Apply statistical methods, datamining & ml algos. to uncover patterns.
- 3) Data visualizaⁿ.

> Benefits

- 1) Data-driven decision making
- 2) Operational Efficiency - identify areas for improvement in processes.
- 3) Risk management.
- 4) Improved customer understanding.

> Application

- 1) Marketing
- 2) Sales
- 3) Finance
- 4) Human Resource.

* ERP & BI - ERP & BI are two critical systems that work together to streamline operations, improve decision making and gain insights.

- > ERP System acts as a central hub for integrating and managing core business process like accounting, inventory, hr.
- > A single database that ensures consistency and accurate accuracy of data across the organization
- > Provides real time access to data, improving visibility & decision making.

> BI refers to technologies, applications and practices for the collection, integration, analysis and presentation of info.

> Integrating ERP systems with BI tools enhances the value of both by combining operational data

> Benefits:

- 1) Enhanced decision making
- 2) Improved performance
- 3) Comprehensive insight
- 4) Forecasting & planning

* BI applications in CRM -

Customer Relationship Management (CRM) systems

Focus on managing all aspects of a customer's relationship with a business, while BI provides tools to extract valuable insights

Application -

- 1) Customer Segmentation - BI analyzes customer data to segment customers into distinct groups based on demographics, purchase behaviour, interests, etc.
- 2) Customer churn - Analyzes customer data and interaction pattern. helps to prioritize customer and personalize marketing strategies.
- 3) Optimizing customer service -
- 4) Decision making .

A) BI application in logistics & production

In the logistics & production sector, a significant volume of data is generated, derived from multiple operations that are carried out everyday.

> BI application -

- 1) Inventory management - Monitor & manage inventory levels, ensuring optimal stock levels are maintained.
Tools - dashboards, predictive analytics.
- 2) Demand Forecasting - To predict future product demand based on historical data.
BI Tools - Time Series analysis, machine learning, predictive modeling.
- 3) Production planning & Scheduling - BI helps analyze historical data on production time, machine performance & material requirements to optimize production schedules.
- 4) Quality Control - Analyze data on product defect and quality control process.
- 5) Transportation management - BI can analyze data on transportation cost, carrier performance, delivery routes. This analysis allows us to find cost saving opportunities, optimizing routes & improved delivery.
- 6) Enhanced customer service.

* Role of BI in Finance

- 1) Financial reporting
- 2) Budgeting & Forecasting
- 3) Risk management.
- 4) Performance management.
- 5) Investment analysis
- 6) Fraud detection & prevention.

* Role of BI in Marketing

- 1) Customer segmentation -
- 2) Customer behavior analysis
- 3) Sales performance management
- 4) Market trend analysis
- 5) Social media analytics.
- 6) Customer churn.

* Role of BI in Telecommunications -

- 1) Customer churn prediction
- 2) Network optimization
- 3) Fraud detection.
- 4) Service usage
- 5) Predictive maintenance
- 6) Network investment & capacity planning

* BI in Salesforce management.

- 1) Sales performance analysis
- 2) Pipeline & forecast management
manage sales pipeline & forecast future sale.
- 3) CRM integration
- 4) Sales activity monitoring
- 5) Churn prediction -
- 6) Customer segmentation.

* BI in HR management

- 1) Recruitment & Talent Acquisition - To streamline hiring process & identifying best candidates.
- 2) Employee performance management.
- 3) Workforce Analytics.
- 4) Employee retention.
- 5) Turnover Analysis.
- 6) Training & Development.
- 7) Workforce planning & optimization.

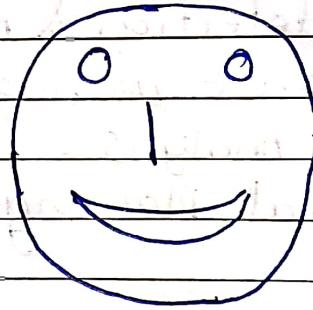
* BI application in Banking

- 1) CRM
- 2) Risk management.
- 3) Credit scoring & loan management.
- 4) Fraud detection & prevention.
- 5) Performance management.
- 6) Transaction Analysis.
- 7) Financial Planning & Analysis.

So This ends here last subject & last exam
of my engineering student life, It was
an incredible journey.
I wish you All the best for your future ☺

Ae

Engineering ke ctrl + c aur ctrl + v kedino ko
ajnida kehete hue ab 'Save as' karke naye chapters
Shuru Karte hain! Farewell coding buddies aur
meri cyber sweetheart (crush), ab 'Run'
Karte hain apna 'life.exe' All the best aur
yaad rakhna, life ke bugs ka fix toh hum sab
kar he lenge.



Rishabh