

# Management Support System

## 1. K-means clustering algorithm

K-means clustering algorithm groups similar objects or data points into clusters. This approach is popular for multivariate numeric data.

- Objective :- i) Grouping similar Data points  
K-Means is designed to cluster data points that share common traits, allowing patterns or trends to emerge.

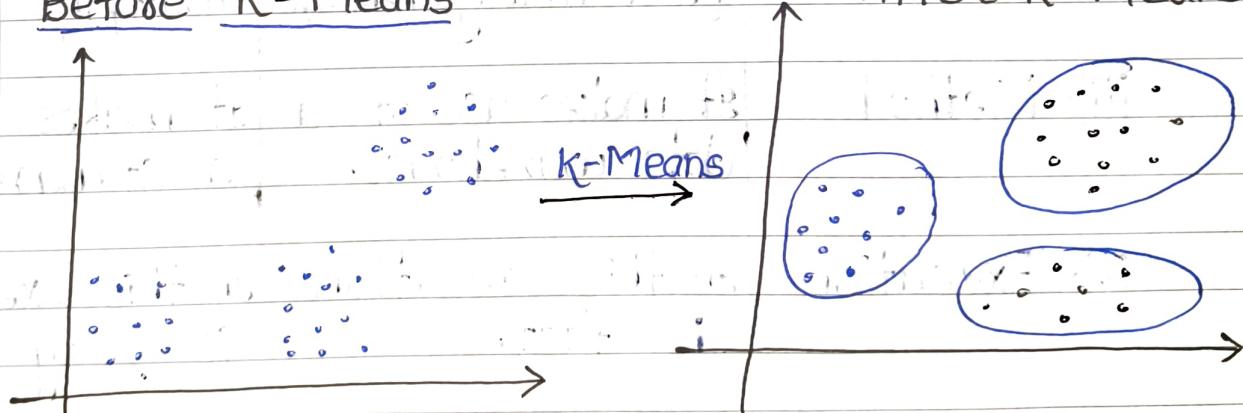
### ii) Minimizing within-cluster Distance

Another objective is to keep data points in each group as close to the cluster's centroid as possible.

Algorithm :- It is an iterative algorithm that divides that unlabeled dataset into k different clusters in such a way that each dataset belongs only one group that has similar properties.

### Before K-Means

### After K-Means



Step 1 :- Select the number k to decide the number of clusters.

Step 2 :- Select random k points or centroid.

Step 3 :- Assign each data point to their closest Centroid, which will form the predefined k clusters.

Step 4 :- Calculate the variance and place a new Centroid of each cluster.

Step 5 :- The model is ready.

## 2. OLAP and OLTP :-

↓ ↳ Online transaction Processing  
Online analytical Processing

Differences :-

Data Source	OLAP	OLTP
i) Data Source	Transcation data based.	Data warehouse or data mart.
ii) Definition	It is well-known as a online database query management Sys.	It is well known as an online database modifying System.
iii) Method	It makes use of a data warehouse.	It makes use of a standard.
iv) Execution	In this execution is fast.	In this execution is very slow.
v) Database design	design with a focus on the Subject.	designed that is focused on the application.

3. What is a cube and what do drill up, roll down and slice and dice mean?

Ans :- The main operational structure of OLAP is based on a concept called Cube. A Cube in OLAP is multi dimensional data structure which might be actual or virtual that allows fast analysis of data. It can also be defined as the capability of efficiently analyse data from multiple perspectives.

- Slice :- A slice is a subset of the cubes corresponding to a single value for one or more members of the dimensions.
- Dice :- The dice operation describes a Subcube by operating a selection on two or more dimension.
- Drill down :- Drilled down or up is a specific OLAP technique where by the user navigates among levels of data record from most summarise to the most detail.
- Roll-up :- A roll-up involves computing all of the data relationship for one or more ~~oranges~~ dimensions to do this a Configurational relationship or a formula might be defined.
- Pivot :- A Pivot is means of changing the dimensional orientation report or ad-hoc query page display.

4. **Apriori Algorithm** :- The most usually Commonly discovered association rule by recursively identify frequent items sets.

Trans no.	SK Us	item sets	Support	item sets	Support	item sets	Support
1	1,2,3,4	1	3	1,2	3	1,2,4	3
1	2,3,4	2	6	1,3	2	2,3,4	3
1	2,3	3	4	1,4	3		
1	1,2,4	4	5	2,3	4		
1	1,2,3,4			2,4	5		
1	2,4			3,4	3		

→ draw this table

→ describes this tables

↳ given a set of example of items to above tables the algorithms attempts to find Subsets that are common to atleast a minimum no. of the item set.

- Apriori uses a bottom up approach where frequent Subsets are extended one item at a time and groups of items at each level are tested against the data for minimum Support. The algorithm terminates number of further successful extension are found.

\* Components of apriori algorithms -

- Support
- Confidence
- Lift

Advantage :- It is used to calculate large itemsets.

Simple to understand and apply.

- Disadvantage :- It is an expensive method.

## 5. Differences - cluster analysis and classification

Ans :- Parameter	Classification	Clustering
Type	used for supervised learning.	used for unsupervised learning.
Basic	Process of classifying grouping the instances the input instances based on their similarity based on their class without the help of labels.	without the help of class.
Complexity	more complex as compared to clustering.	less complex to compared to classification.
Methods	decision tree, logistic regression	k-mean clustering, hierarchical clustering.
Data Requirement	works with unlabeled data.	Requires labeled data for training.
Output	Produces clusters without specific labels.	Produces discrete known labels for each data point.

## 6. General algorithm of decision tree with example.

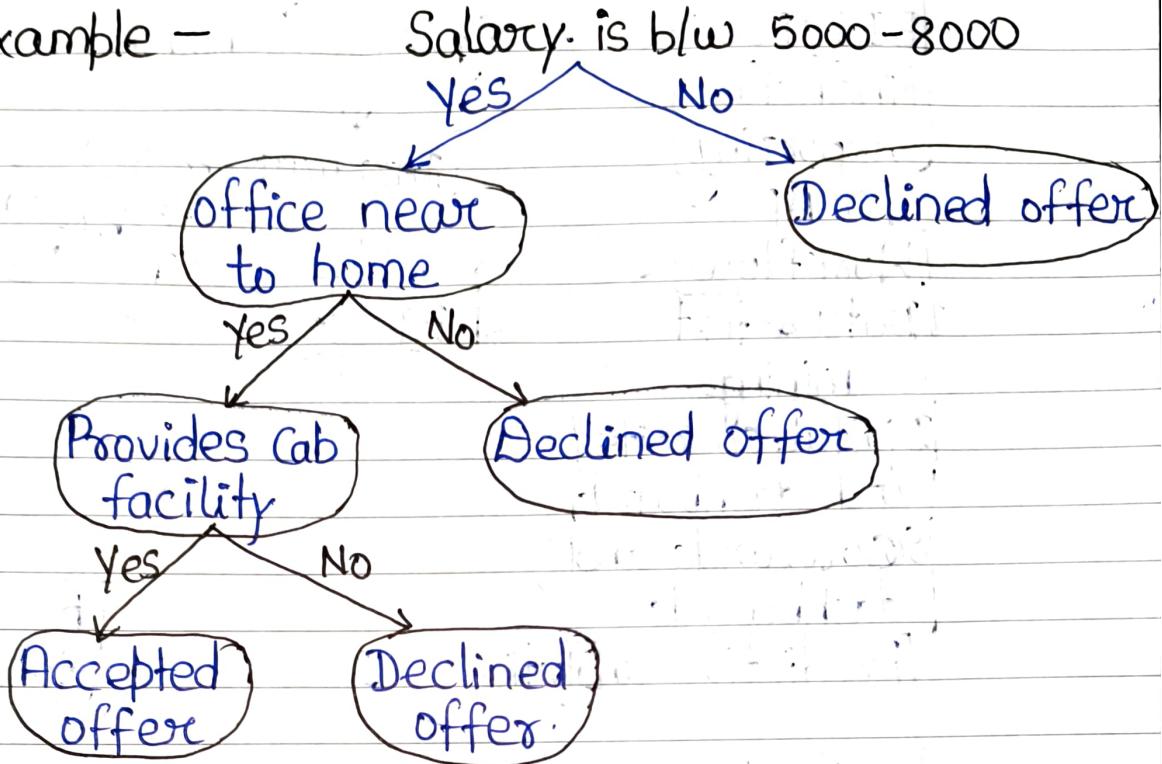
Ans : A decision tree is a type of Supervised Learning algorithm that is commonly used in machine learning to model and predict outcomes based on input data. It is a tree-like structure where an internal node test on attribute, each branch corresponds to attribute value and each leaf node represents the final decision or prediction. The decision tree algorithm falls under the category of supervised learning. They can be used to solve both regression and classification.

- Root node
- internal nodes
- leaf nodes
- Branches
- splitting
- Parent node
- child node

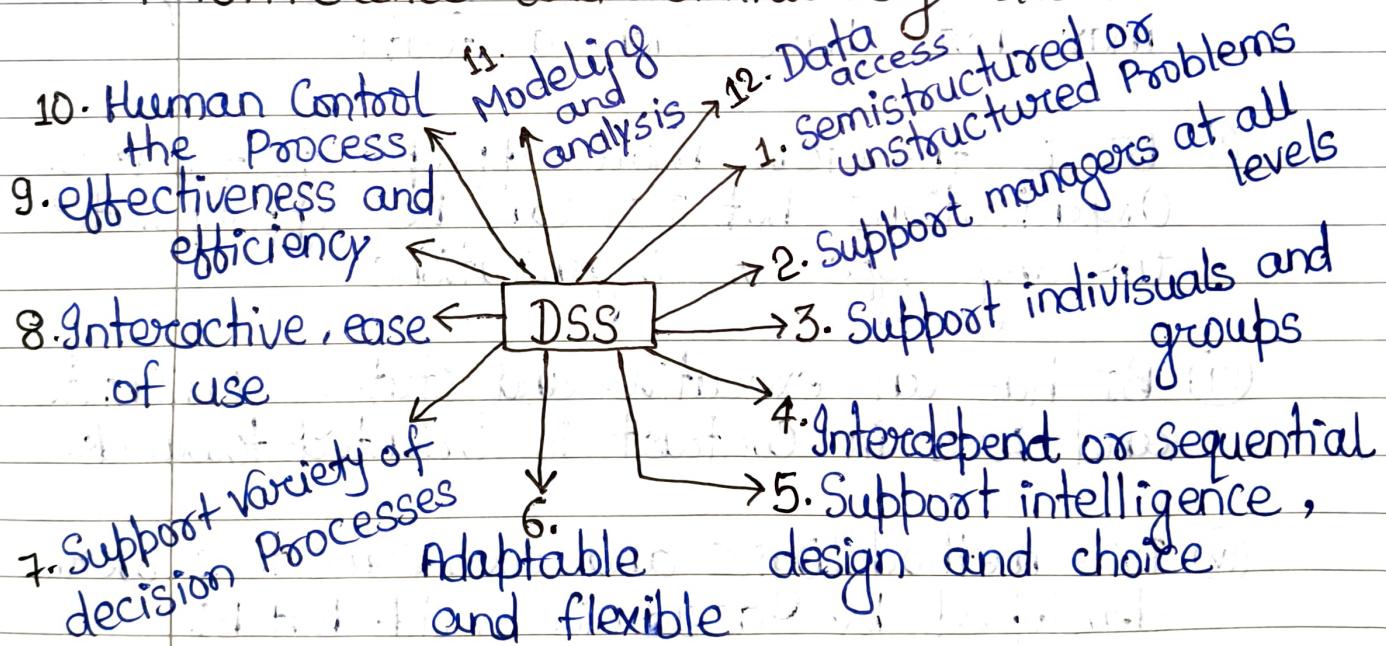
Algorithm :- ① Begin the tree with the root node →  
Says X which contains the complete datasets.

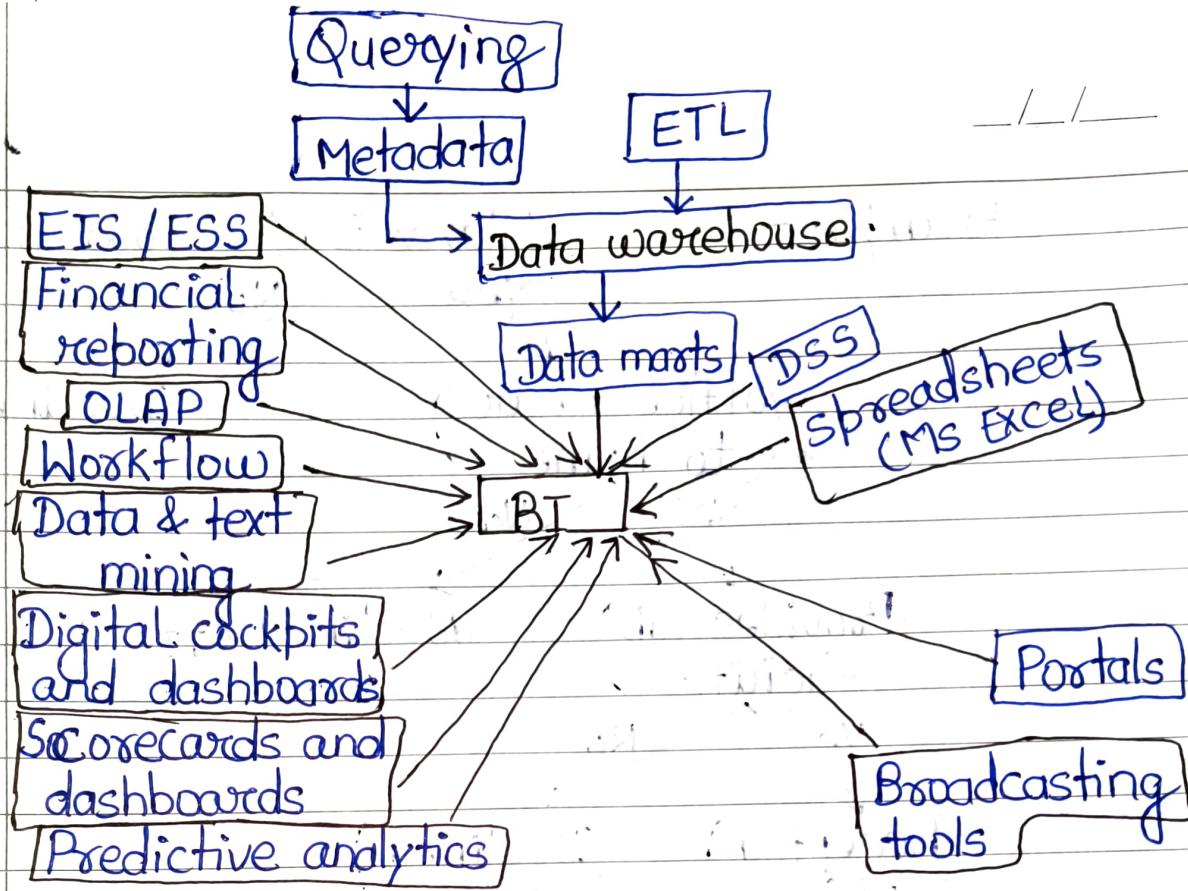
- ① find the best attribute in the dataset using attribute selection.
- ② Divide the S into subsets that contains possible values for the best attributes.
- ③ Generate the decision tree node, which contains the best attribute.

Example -



## 7. Difference and Similarity b/w DSS & BI.





### Similarities and differences -

- (i) Their architecture are very similar because BI evolved from DSS. but DSS may or may not have such a feature.
- (ii) BI is therefore more appropriate for large Organizations but DSS can be appropriate to any type of Organizations.
- (iii) BI has an executive and strategy Orientation and DSS in contrast is Oriented towards.
- (iv) DSS methodologies and even some tools were developed by academic world. But BI methodologies and tools were developed by s/w Companies.

## 8. Gini index and how does it measure?

Ans:- Gini index is a metric to measure how often a randomly chosen element would be incorrectly identified.

- It means an attribute with a lower gini index should be preferred.
- Sklearn supports "Gini" criteria for Gini index and by default it takes "gini" value.

$$\text{Gini}(S) = 1 - \sum_{i=1}^c p_i^2$$

The Gini index is a measure of the inequality or impurity of a distribution. Commonly used in decision trees and other machine learning algorithms. It ranges from 0 to 0.5 where 0 indicates a pure set and 0.5 indicates a maximally impure set.

- In decision trees, the Gini index is used to evaluate the quality of a split by measuring the difference b/w the impurity of the Parent node and the weighted impurity of child node.

## 9. What an Organization should Consider before making a decision to purchase data mining Software?

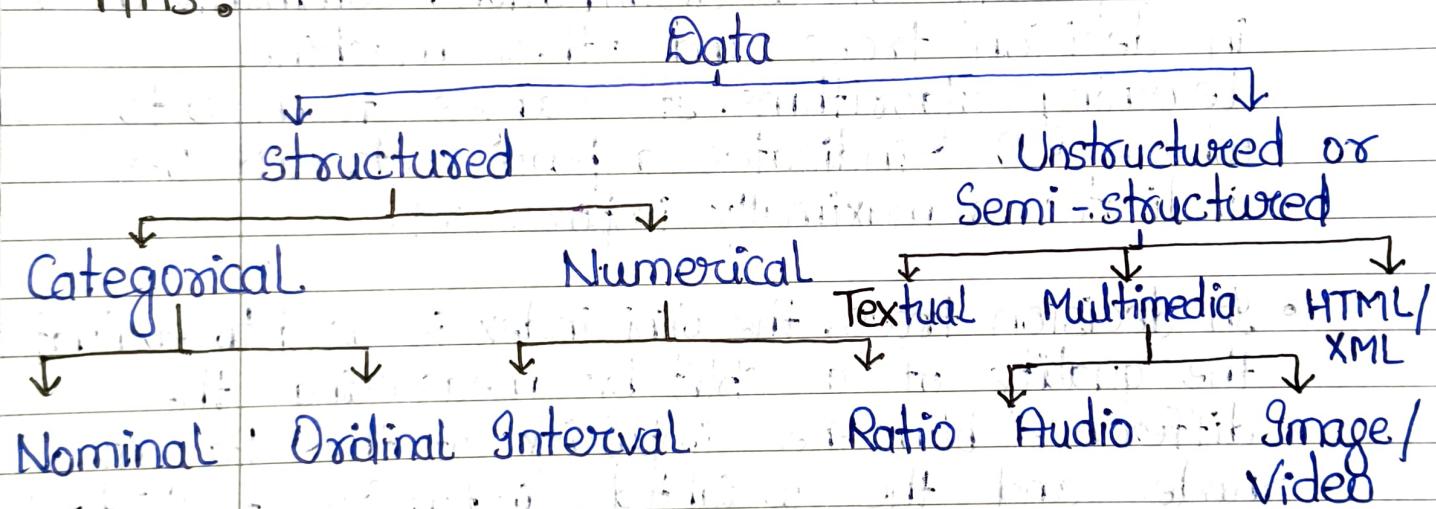
Ans: Before Purchasing data mining Software, an Organization should consider things like :-

- Data quality and accuracy :- How well the S/w can handle the quality and size of the data.

- Usability :- How easy the SW is to use.
- Visualization :- what visualization features the SW offers.
- Cost :- How costly the SW will be to invest.
- Efficiency and Scalability :- How efficient and Scalable the SW is.
- Support Services :- what Support Services are available for the Software.
- Vendor Reputation :- The reputation of the SW vendor should also be considered.

## 10. Taxonomy of Data mining

Ans:-



- Categorical data represent the labels of multiple classes used to divide a variable into specific groups.
- Nominal data contain measurements of simple codes assigned to objects as labels, which are not measurements.

- Ordinal data contains codes assigned to objects or events as labels that also represent the rank order among them.
- Numeric data represent the numeric values of specific variables. e.g - numeric value
- Interval data are variables that can be measured on interval scales.
- Ratio data includes measurement variables commonly found in the physical sciences and engineering. e.g - Mass, length, time  
other datatypes including textual, spatial, imagery, and voice, need to be converted into some form of categorical or numeric representation before they can be processed by data mining algorithms.  
Data can also be classified as static or dynamic.

The nominal or ordinal variables are converted into numeric representations using some type of 1-of-N Pseudo variables.

## 10. How the DSS are Supported?

Ans:- A Decision Support System is an interactive Computer-based System that assists individuals or organizations in making informed decisions by analyzing large volumes of data. Provides insights and helping predict Outcomes.

Supported? :- Here we relate specific technologies to the decision making Process. Databases, data marts, and especially data warehouses are important technologies in Supporting all phases of decision making.

### ① • Support for the intelligence Phase

The Primary requirement of decisions Support for the intelligence Phase is the ability to Scan external and internal information Sources for opportunities and Problems and to interpret what the Scanning discovers.

- Decision Support / BI technologies Can be very helpful. Graphical Systems information System can be utilized either as stand-alone Systems or integrated with these Systems so that a decision maker can determine opportunities and Problems in a spatial Sense. Another aspect of identifying internal Problems and Capabilities involves monitoring the Current status of operations.

Expert Systems in Contrast can render advice regarding the nature of a Problem, its classification, its seriousness, and the like.

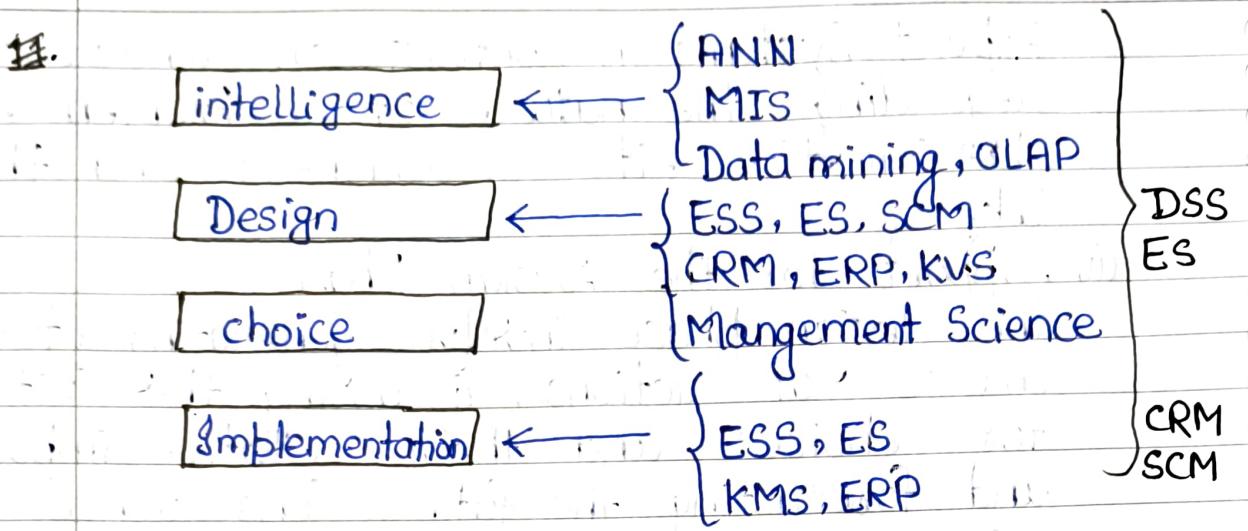
### (ii) Support for the Design Phase

The design phase involves generating alternatives Course of actions, discussing the Criteria for choices and their relative importance and forecasting the future Consequences of using various alternatives. Several of these activities can use standard models provided by a DSS.

### (iii) choice Phase :- In addition to Providing models that rapidly identify a best or good enough alternatives, a DSS can support the choice Phase through what-if and goal-seeking analyses.

### (iv) implementation Phase :- DSS can be used in implementation activities Such as decision Communication, explanation, and Justification. Implementation Phase DSS benefits are partly due to the vividness and detail of analyses and reports. Reporting Systems and others tool variously labeled as BAM, BPM, KMS, EIS, ERP, CRM and SCM are all useful in tracking how well an implementation is working.

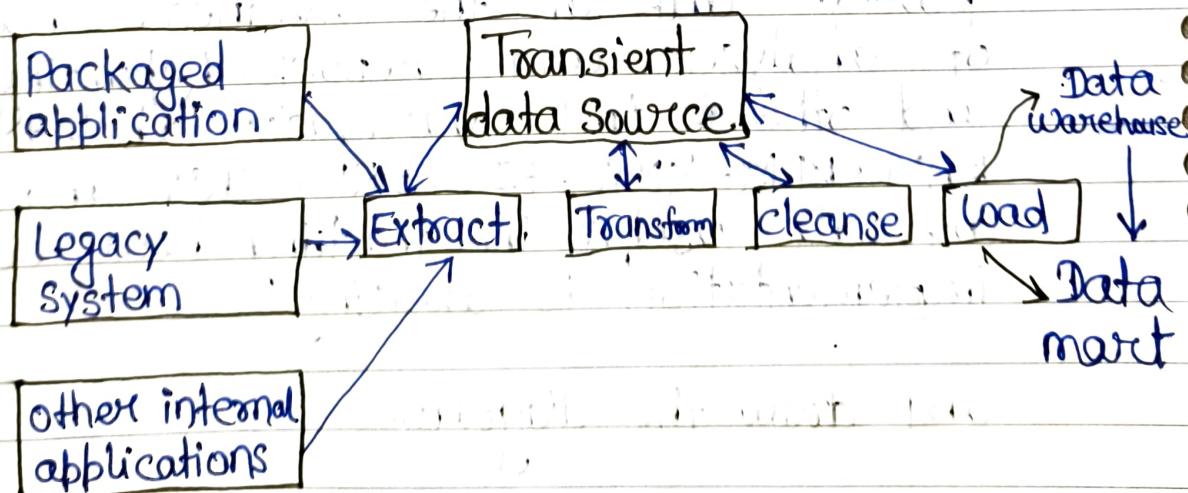
- and draw the diagram.



## 11. Data Integration and the Extraction, Transformation and Load (ETL) Process

Ans :- Data integration Comprises three major Processes that, when correctly implemented, Permit data to be accessed and made accessible to an array of ETL and analysis tools and the data warehousing environment: data access , data federation and change capture .

- Enterprise application integration (EAI)
- Service - Oriented architecture (SOA)
- Enterprise information integration (EII)
- Extraction , transformation , and load (ETL)



At the heart of the technical side of the data warehousing Process is extraction, transformation and load (ETL).

Extraction ETL technologies, which have existed for some time, are instrumental in the Process and use of data warehouses. The ETL Processes is an integral Component in any data-centric Project. IT managers are often faced with challenges because the ETL Process typically Consumes 70% of the time in a data-centric Project. IT managers are often faced with challenges because the ETL Process typically Consumes Transformation Occurs by using rules or lookup tables or by Combining the data with other data.

- Data transformation tools are expensive.
- Data transformation tools may have a long learning Curve.

\*\*

Data mining

→ Prediction

→ classification

→ Regression

→ Association

→ link analysis

→ Sequence analysis

→ clustering

→ Outlier analysis

Learning method → Popular algorithm

i) Supervised → classification and regression  
trees, ANN, SVM, GA

ii) Supervised → Decision trees, ANN / MLP, SVM,  
Rough sets, GA

Supervised → linear / Nonlinear Regression,  
Regression trees, ANN / MLP, SVM

iii) Unsupervised → Apriori, OneR, ZeroR, Eclat  
→ Expectation maximization, Apriori algorithm

iv) Unsupervised → k-Means ANN / SOM  
→ k-means, Expectation maximization.

## Fill in the blanks -

① \_\_\_\_\_ is an important Component of text mining and is soft field of AI and Computational linguistic it study of Problem of understanding natural human language.

Ans:- NLP

② \_\_\_\_\_ is One or more web pages that provides a collection of links authority Pages.

Ans:- Hub Page

③ \_\_\_\_\_ mining is the Process of extracting useful information from the links ambedded in web documents.

Ans:- Web structure

④ \_\_\_\_\_ Segmentation is one of the unsupervised machine learning methods where in data is put into Several groups based on their similarity.

Ans:- clustering

⑤ \_\_\_\_\_ deals with the kind of uncertainty and partial information that is human in nature.

Ans:- fuzzy logic

⑥ \_\_\_\_\_ is a relational approach to design and developed of a model based System.

Ans - Relational model based System

vii) if the simulation results do not match the equitance and Judgement of decision maker than a confidence can occur in the result.

viii) Continuous distributions are situations with limited no. of possible events that follows density function Such as the normal distributions.