

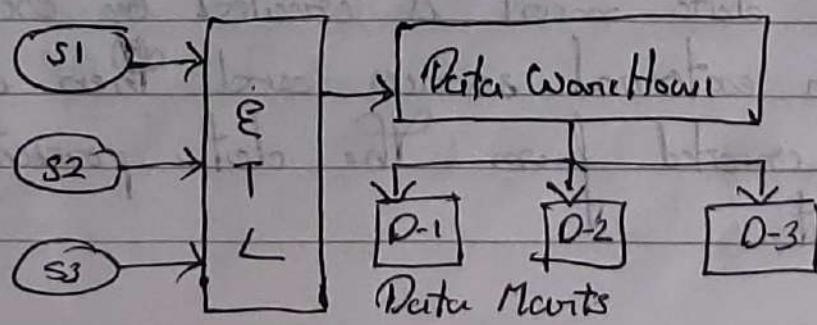
## Data Mart:-

- A data mart is a specialized subset of a data warehouse focused on a specific functional area or department within an organization.
- It provides a simplified and targeted view of data addressing specific reporting and analytical needs.
- They are organized around specific subjects, such as sales, customer data, or product information and are structured, transformed and optimized for efficient querying and analysis within the domain.

## -Types of Data Mart:-

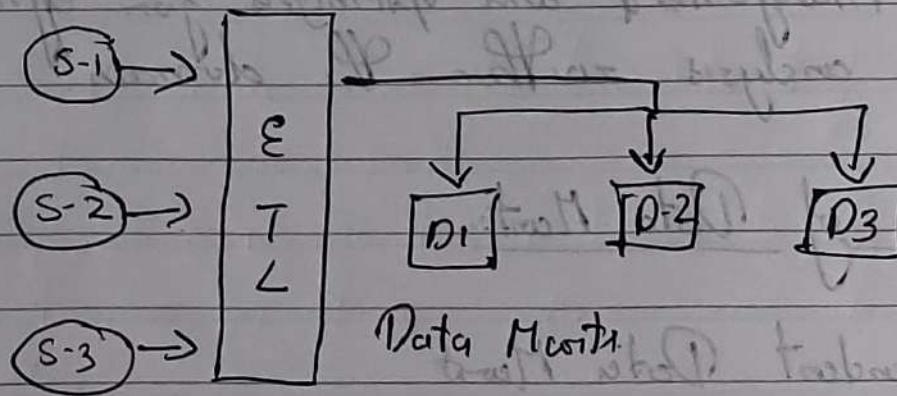
- i. Dependent Data Mart
- ii. Independent Data Mart
- iii. Hybrid Data Mart

## Dependent Data Mart..



- Created by extracting data from central repository i.e Datawarehouse
- First data warehouse is created by extracting data (through ETL Tool) from external sources and then data mart is created from data warehouse.
- It uses the top-down approach.
- Used by big organizations

### Independent Data Mart



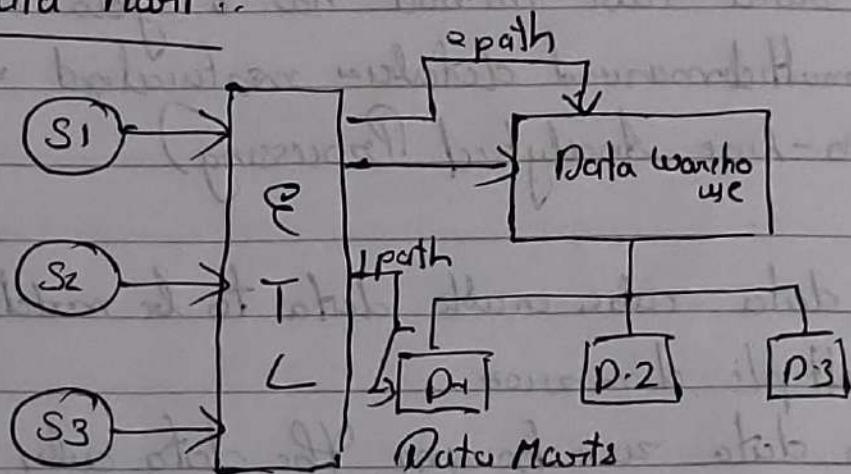
Created directly from external sources instead of data warehouses.

- First data mart is created by extracting data from external sources and then data warehouse is created from the data present in data mart.

→ Used Bottom up approach.

→ Used by Smaller organizations

### Hybrid Data Mart:



- Created by extracting data from operational source or from data warehouse
- Path 1 reflects accessing data directly from external sources and path 2 reflects dependent data model of data mart.

### Multidimensional data Modelling

If is a data modeling technique used in data warehouses to organize data in the database in an efficient manner to analyze future trends and patterns.

## Data Cube:-

- When data is grouped or combined in multidimensional matrix, it is called Data Cube.
- The data cube method has a few alternative names as multidimensional database, materialized view and OLAP (On-line Analytical Processing).
- A data cube enable data to be modeled and viewed in multiple dimensions.
- In data warehousing, the data cubes are n-dimensional.
- The cuboid which holds the lowest level of summarization is called a base cuboid.
- The topmost 0-D cuboid, which holds the highest level of summarization, is known as the apex cuboid.
- In this example, this is the total sales in dollars sold, summarized over all four dimension.
- Data cube is represented by dimensions and facts.
- Each dimension can have a table related to it and is known as a dimension table.
- Eg.: Dimension Table for an item can include the attributes items name, Brand and type.

- > A multidimensional data model is generally organized around a central design called facts
- > The fact table includes the names of the facts or measures and keys to each of the associated dimension tables.

## Fact Table - Sales

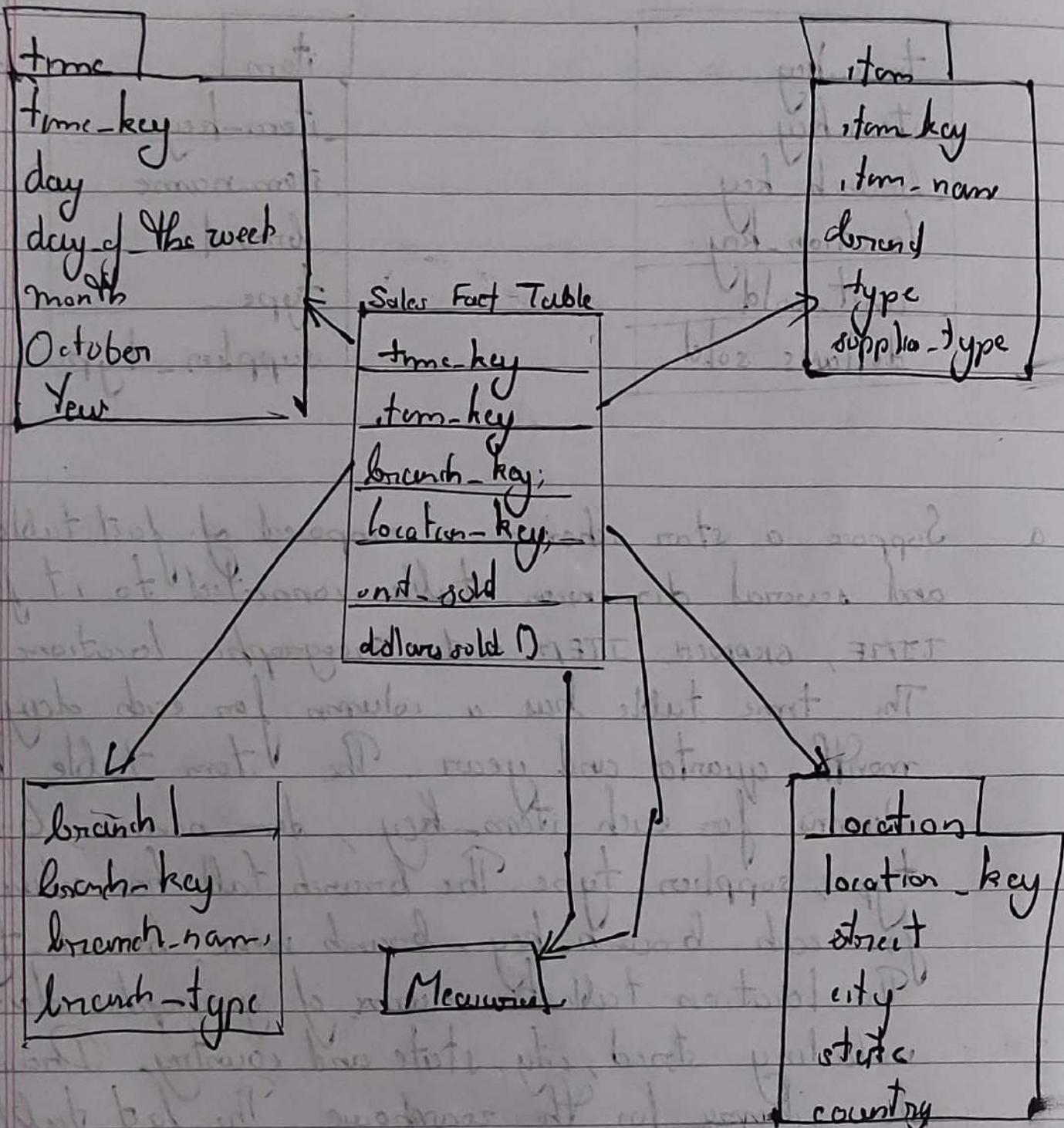
time-key
item-key
branch-key
location-key
units-sold
dollars-sold

## Dimension Table - Item

item
item-key
item-name
brand
type
supplier-type

a Suppose a star schema is composed of fact table, sales, and several dimensions tables connected to it for TIME, BRANCH, ITEM and geographic locations.

The time table has a column for each day, month, quarter and year. The item table has columns for each item-key, item-name, brand, type, supplier-type. The branch table has columns for each branch-key, branch-name, branch-type. The location table has columns of geographic data including street, city, state and country. Draw the star schema for the warehouse. The fact table contains two measures units-sold and dollars-sold.



- a) A data warehouse for a shop consists of the following 3 dimensions. Store, time, product with 2 measure units and price. Time dimension (year, quarter, month), store dimension store name, city, state, region, Product dimension (class, type)
- Draw the star schema for the warehouse.

Star Schema

Store Dimension

Fact Table

Store-Key

Product-Key

Period-Key

Units

Price

Time Dimension

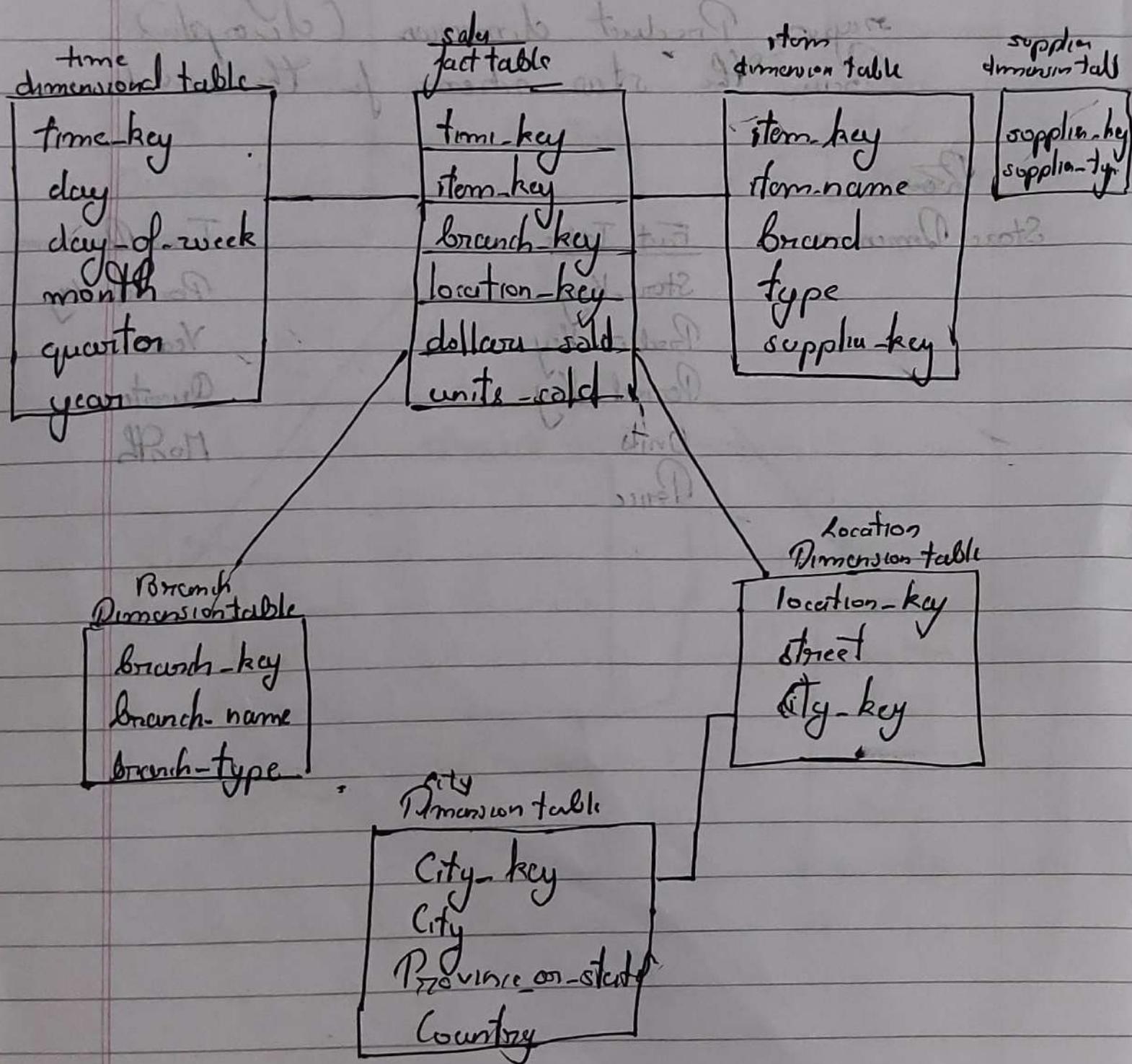
Period-Key

Year

Quarter

Month

- \* Snowflake Schema has only one fact table and multiple dimension tables. The dimension tables are connected to sub-dimension tables if required.



- Fact Constellation has multiple fact tables with dimension tables and the dimension tables share the same packet.

Time

gallu

10m

shopping

## dimension table

justable

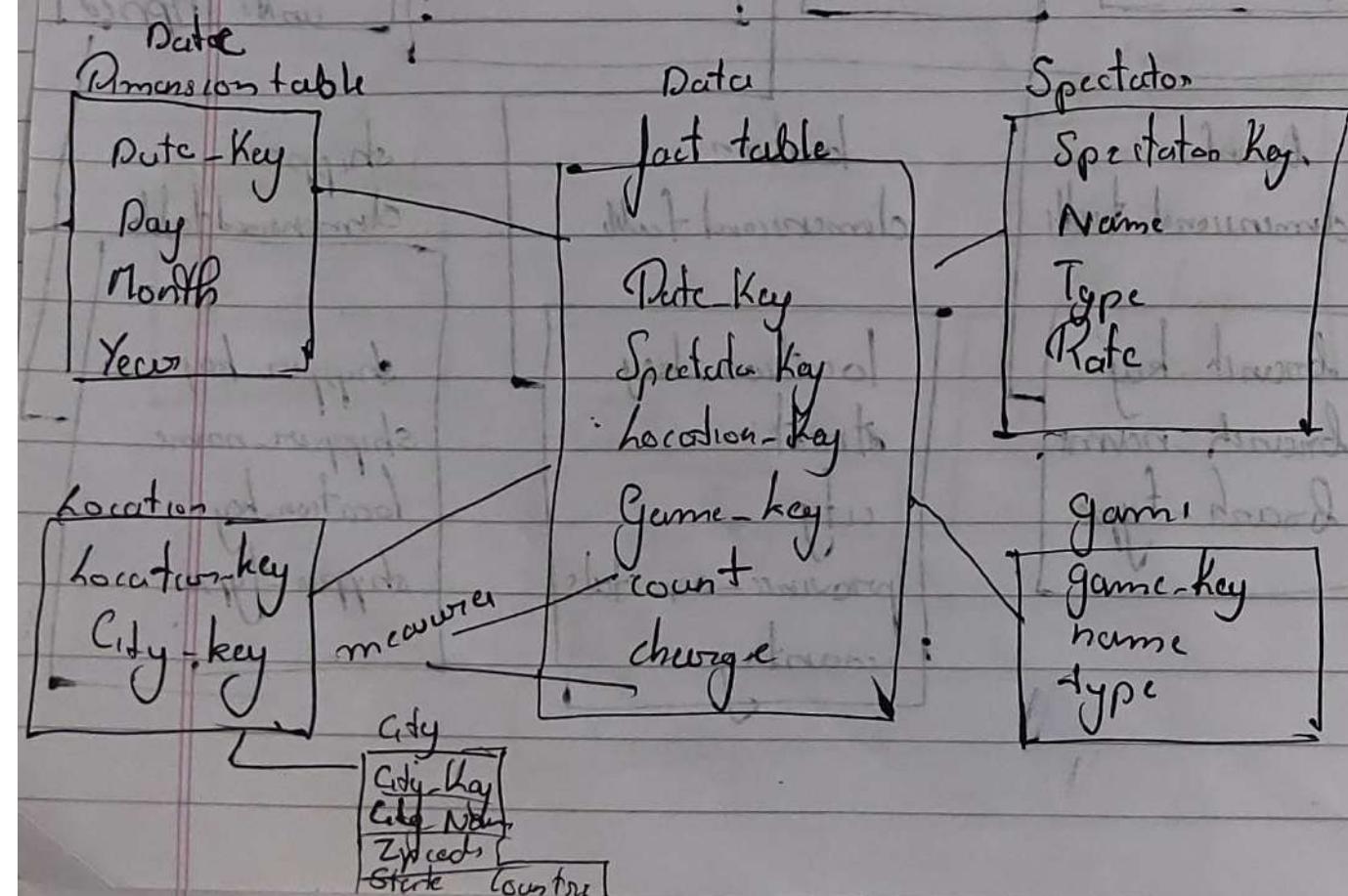
## dimensions of cable

fact falls



a) Suppose a data warehouse consists of 4 dimensions date, spectator, location and game with two measures count and charge where charge is the fee that a spectator plays when watching a game on a given day. Spectators may be students, adults or seniors with each category having its own charge rate. At higher conceptual level, location stores the details of city, city.name, zipcode, state and country. Identify a suitable multidimensional model for the given data warehouse.

It is snowflake ∵ location has sub dimensions



## OLAP

- OLAP stands for online analytical processing
- In the multidimensional model, the records are organized into various dimensions, and each dimension includes multiple levels of abstraction described by concept hierarchy.
- It supports users with the flexibility to view data from various perspectives.

## OLAP Operations

- i Roll Up:-
- ii Drill Down
- iii Slice and Dice
- iv Pivot

### Roll Up

- Roll Up is also known as consolidation or aggregation and it can be performed in two ways:
  - ↳ Reducing Dimensions
  - ↳ Climbing Up Concept Hierarchy
- Concept Hierarchy is a system of grouping things

Based on their order or level.

Eg:- Roll up on location from cities to countries

→ New Jersey and Los Angeles are rolled up into country USA.

## Module 3

### Machine learning

Machine learning is a growing technology that enables computer to learn automatically from past data.

It uses various algorithms for building mathematical models and making prediction using historical data or information.

The term machine learning was introduced by Arthur Samuel in 1959.

Def → Machine learning enables a machine to automatically learn from data improve performance from experiences and predict things.

Type of ML

- i) Supervised learning
- ii) Unsupervised learning
- iii) Reinforcement learning

## Supervised learning:

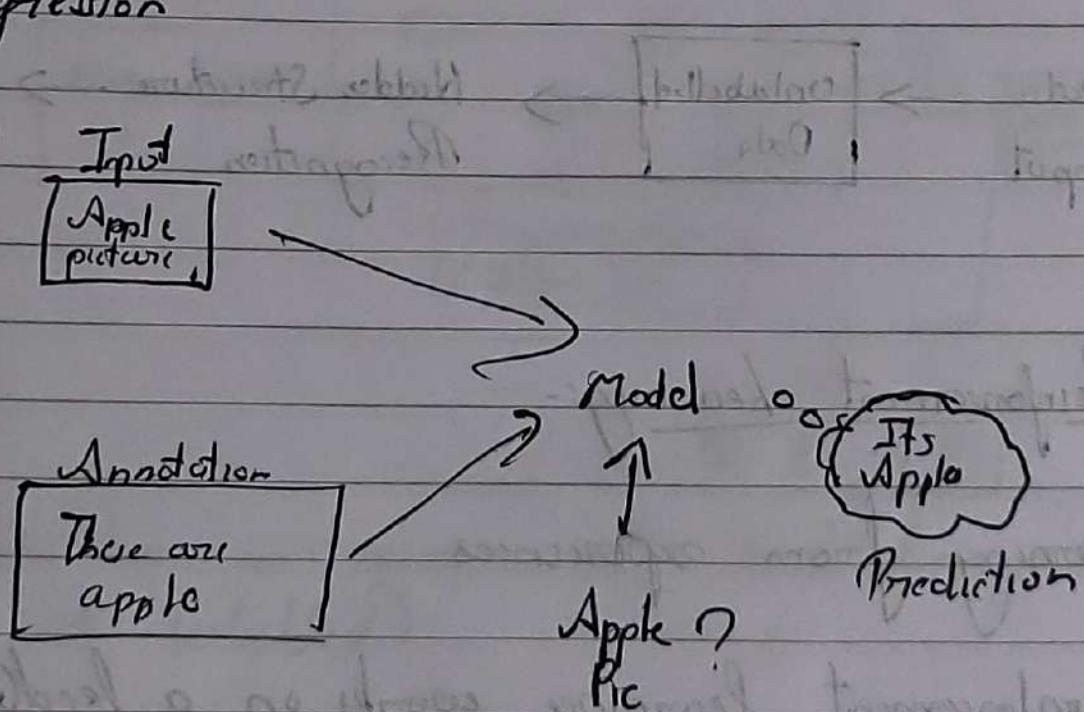
Supervised machine learning is based on supervision.

In this technique, we train the machine using the "labelled" dataset and based on the training, machine predicts the outcome.

The main goal of the supervised learning technique is to map the input variable ( $x$ ) with output variable  $y$ .

### Supervised Learning Types:

- i Classification
- ii Regression

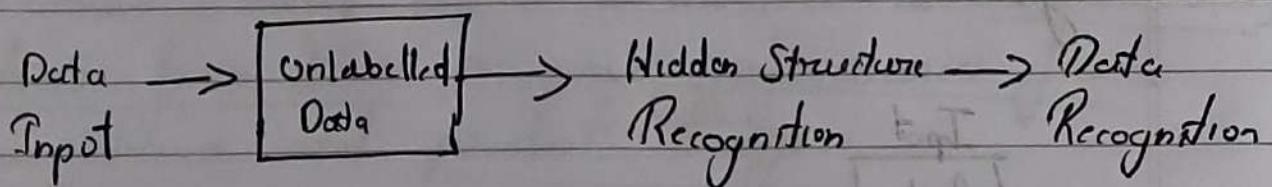


## Unsupervised Learning:-

- There is no supervision
- The machine is trained using unlabelled dataset and the machine predicts the output.
- Aim of this algorithm is used to group or categorize the dataset according to similarities, pattern and differences.

Two Type

- i) Clustering
- ii) Association



## Reinforcement Learning:-

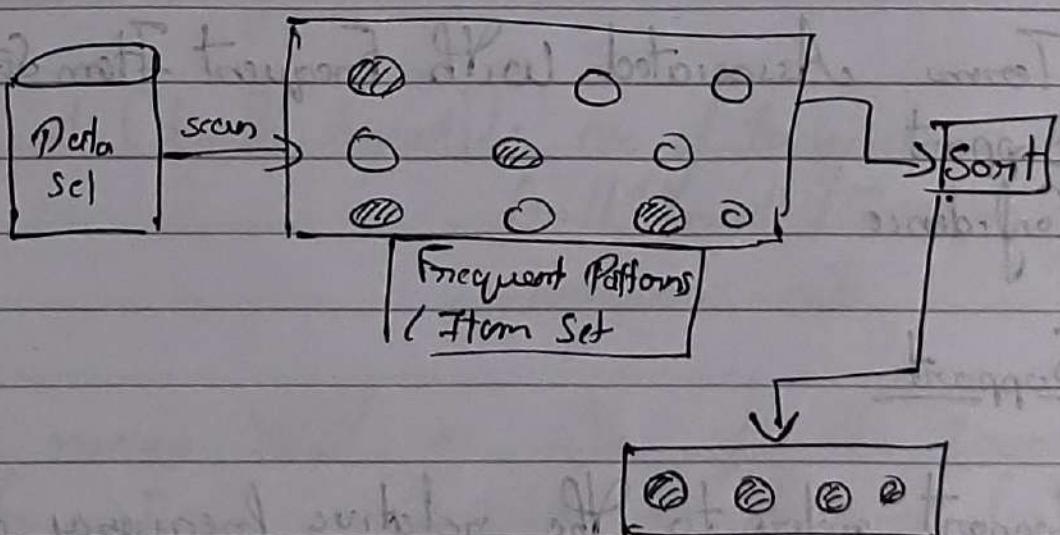
- learning from experiences
- Reinforcement learning works on a feedback based process in which an AI agent (software)

component) automatically explore its surroundings, learning from experiences and improving its performance.

- The agent gets rewarded for each good action and get punished for each bad action, hence the goal of reinforcement learning agent is to maximize the rewards.

### Frequent Pattern Mining :-

Frequent pattern refers to itemset, subsequences or substrings that appear frequently in a dataset.



Section (This is an unsupervised dataset because the data is not labelled)

## Applications

### → Market Basket Analysis:

- In DM, it is to analyze the combination of products which has been bought together.
- It is used by retailers to increase sales by better understanding customer purchasing patterns.

It involves analyzing large data sets, such as purchase history to reveal product groupings, as well as products that are likely to be purchased together.

- Terms Associated with Frequent Item Set
- Support
- Confidence

### Support

- Support refers to the relative frequency of an item in a dataset.
- The support of an itemset is the number of transactions in which the itemset appears, divided

by the total number of transactions

Support is calculated as follows

$$\rightarrow \text{Support}(X) = (\text{Number of transactions containing } X) / \text{Total no of transactions.}$$

where  $X$  is the item set for which we are calculating support.

For example, suppose we have a dataset of 1000 transactions and the itemset {milk, bread} appears in 100 of those transactions.

The support of the itemset {milk, bread} would be calculated as follows

$$\begin{aligned} \text{Support}(\{\text{milk, bread}\}) &= \text{No of transactions containing } \{\text{milk, bread}\} / \text{Total no of transactions} \\ &= 100 / 1000 = 10\% \end{aligned}$$

This means that in 10% of the transactions, the items milk and bread were both purchased

Algorithms that can be used for classification using frequent pattern mining

- i Apriori Algorithm
  - ii FP Growth Algorithm
- a Consider The dataset

Transaction id | items

T <sub>1</sub>	hot dogs, buns, ketchup
T <sub>2</sub>	hot dogs, buns
T <sub>3</sub>	hot dogs, coke, chips
T <sub>4</sub>	chips, coke ..?
T <sub>5</sub>	chips, ketchup
T <sub>6</sub>	hotdog, coke, chips

Find The frequent item set and generate association rules using apriori algorithm. Assume The minimum support threshold is 33.33% and minimum confidence is 60%.

Minimum Support Count

$$\text{No of transactions} \times \frac{\text{support}}{100}$$

$$6^3 \times \frac{33.33}{100} = 1.99 = 2$$

Step 1

<u>Item</u>	<u>Support Count</u>
-------------	----------------------

Hot Dogs	4
Buns	2
Ketchup	2
Coke	3
Chips	4

If any item having the count less than minimum support count, then we need to remove that item (Data pruning)

Step 2

<u>Item</u>	<u>Support Count</u>
-------------	----------------------

Hot Dogs, Buns	2
Hot Dogs, Ketchup	1
Hot Dogs, Coke	2
Hot Dogs, Chips	2
Buns, Ketchup	1
Buns, Coke	0
Buns, Chips	0
Ketchup, Coke	0
Ketchup, Chips	1
Coke, Chips	3

Step 3

Data Mining

<u>Item</u>	<u>Support Count</u>
Hot Dog, buns	2
Hot Dogs, coke	2
Hot Dogs, chips	2
Coke, chips	3

Step 4

Combination of 3

<u>Item</u>	<u>Support Count</u>
Hot Dogs, Buns, Ketchup	1
Hot Dogs, Buns, Coke	0
Hot Dogs, Buns, Chips	0
Hot Dogs, Ketchup, Coke	0
Hot Dogs, Ketchup, Chips	0
Hot Dogs, Ketchup, Eggs	2
Buns, Ketchup, Coke	0
Buns, Ketchup, Chips	0
Ketchup, Coke, Chips	0
Bun, Coke, Chips	0

Step 5

Data Mining

Item                          Min Support

Hotdogs, chips, coke      2

Frequent itemset: {hotdogs, chips, coke}

### Association Rule

$$\{\text{hotdogs, coke}\} \Rightarrow \{\text{chips}\}$$

$$\text{Confidence} = \frac{\text{Support count}(\{\text{hotdogs, chips}\})}{\text{Support count}(\{\text{hotdogs, coke}\})} \times 100$$

$$= \frac{2}{2} \times 100 = 100\%$$

$$\{\text{hotdogs, chips}\} \Rightarrow \{\text{coke}\}$$

$$\text{Confidence} = \frac{\text{Support count}(\{\text{hotdogs, coke}\})}{\text{Support count}(\{\text{hotdogs, chips}\})} \times 100$$

$$= \frac{2}{2} \times 100 = 100\%$$

$$\{\text{coke, chips}\} \Rightarrow \{\text{hotdogs}\}$$

$$\text{Confidence} = \frac{\text{Support count}(\{\text{coke, hotdogs}\})}{\text{Support count}(\{\text{coke, chips}\})} \times 100$$

$$= \frac{2}{3} \times 100 = 66.66\%$$

$\{ \text{hotdogs} \} \Rightarrow \{ \text{coke, chips} \}$

Confidence:  $\frac{2}{4} \times 100 = 50\%$

$\{ \text{coke} \} \Rightarrow \{ \text{hotdogs, chips} \}$

Confidence =  $\frac{2}{3} \times 100 = 66.66\%$

$\{ \text{chips} \} \Rightarrow \{ \text{hotdogs, coke} \}$

Confidence =  $\frac{2}{4} \times 100 = 50\%$

### Strong Association Rule

- 1  $\{ \text{hotdogs, coke} \} \Rightarrow \{ \text{chips} \}$
- 2  $\{ \text{hotdogs, chips} \} \Rightarrow \{ \text{coke} \}$
- 3  $\{ \text{coke, chips} \} \Rightarrow \{ \text{hotdogs} \}$
- 4  $\{ \text{coke} \} \Rightarrow \{ \text{hotdogs, chips} \}$
- 5  $\{ \text{chips} \} \Rightarrow \{ \text{hotdogs, coke} \}$

## Module IV

### Classification & Prediction

→ There are two forms of data analysis that can be used to extract models describing important classes or predict future data trends.

They are:-

- ↳ Classification
- ↳ Prediction.

These are the two main method used to mine the data.

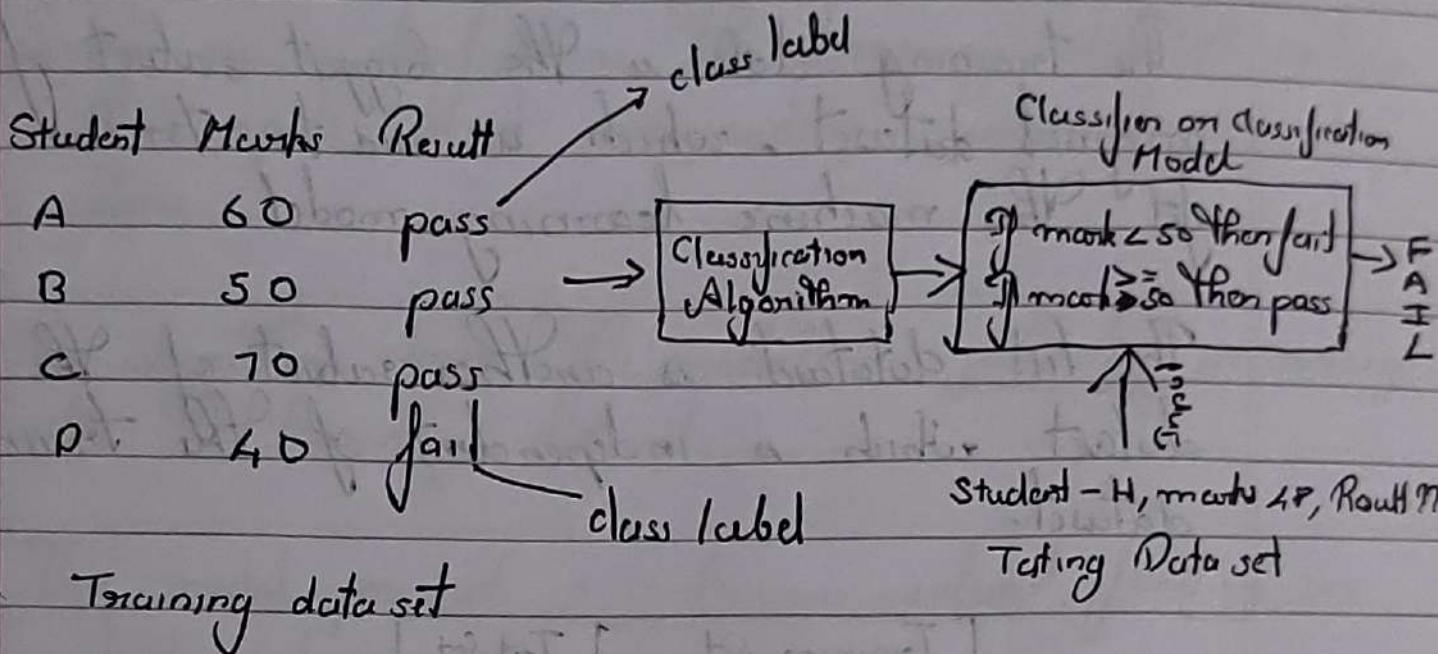
#### Classification:-

It is for converting large datasets to useful info.

- Application → To classify people with or without a certain disease
- To classify spam messages

Classification is the process of dividing the dataset into different categories or groups by

## adding class labels



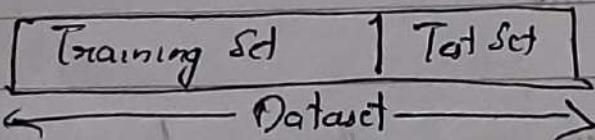
→ Supervised learning always has a training set  
 $\therefore$  Classification is a supervised learning algorithm

- First a set of data is used as training data
- The set of input data and the corresponding output are given to the algorithm
- So the training dataset includes the input data and the associated class labels.
- Using the training dataset, the algorithm derive a model or classifier.

## Training and Testing Data:

The training data is the biggest subset of the original dataset, which is used to train or fit the machine learning model.

The test dataset is another subset of the original dataset which is independent of the training dataset.



The training dataset is generally larger in size compared to testing dataset. General ratios of splitting train and test datasets are 80:20, 70:30, or 90:10.

## Prediction:-

Forecast about a future event.

Eg Rainfall Prediction,

In this method, we need to predict the missing data for a new observation depending

on the previous data

Same as in classification. The training dataset contains the inputs and corresponding numerical output values.

The algorithm derives the model on a predictor according to the training dataset.

The model predicts a continuous-valued function or ordered value.

## Bayesian Classification :- (Naive Bayes Classifier)

- Naive Bayes is a probabilistic machine learning algorithm that can be used as a wide variety of classification tools.
- Bayesian classification <sup>classifier.</sup> can predict the probability that a given talk belongs to a particular class.
- Based on Baye's Theorem.
- More speed and accuracy as compared to decision tree classifier.

## Baye's Theorem

To calculate the probability of an event, based on prior knowledge related to that event.

$$\text{Expressed as} \rightarrow P(A|B) = P(B|A) * P(A) / P(B)$$

$P(A|B)$  is the posterior probability. (Probability of hypothesis A - on the observed event B,

$P(B|A)$  is likelihood probability. (Probability of the evidence given that the probability of a hypothesis is true.

$P(A)$  and  $P(B)$  are the probabilities of event A and B

Q) Consider the dataset

Find out - class label on the instance - outlook: sunny

	outlook	Play
0	Rainy	Yes
1	Sunny	Yes
2	Overcast	Yes
	Overcast	Yes
	Sunny	No