

Data Warehouse

→ Provide tool for business executives to systematically organize, understand and use their data to make strategic decisions.

Key features of data warehouse (Imp. 6x1 = 6)

- 1) Subject oriented: It is organized around major subjects such as customers, suppliers, product and sale.
- 2) Integrated: Data warehouse is constructed by integrating multiple heterogeneous sources such as flatfiles.
- 3) Time variant: Data are stored to provide info. from a historical perspective (past 5 to 10 years).
- 4) Non volatile: Data warehouse is a physically separate store of data transformed from application data found in the operational environment.

Difference b/w operational DB systems and Data warehouse

- OLTP → cover day to day operations. (Online transaction Processing)  
e.g., Banking, inventory, purchasing.
- DW → Data analysis and decision making. (OLAP)

| OLTP Feature             | OLTP                             | OLAP  |
|--------------------------|----------------------------------|---|
| 1) Characteristics:      | Operation processing             | Information processing                                |
| 2) Orientation           | Transaction                      | Analysis  |
| 3) User                  | Clark, DBA, DB professional      | Knowledge worker<br>e.g., manager, executive, analyst |
| 4) Function              | Day to day operations            | Long term information requirements.                   |
| 5) DB design             | ER based<br>↓<br>Entity-Relation | star or snowflake,<br>subject oriented                |
| 6) Access                | Read and Write                   | Mostly Read   |
| 7) No. of records stored | 10 (thrs)                        | Millions  |

- 8) No. of users 1000  
 9) DB size 100 MB to GB  
 10) Data storage stores current data
- 100  
 100 GB to TB  
 stores historical data
- 02/11/2023

### Multidimensional Data Model

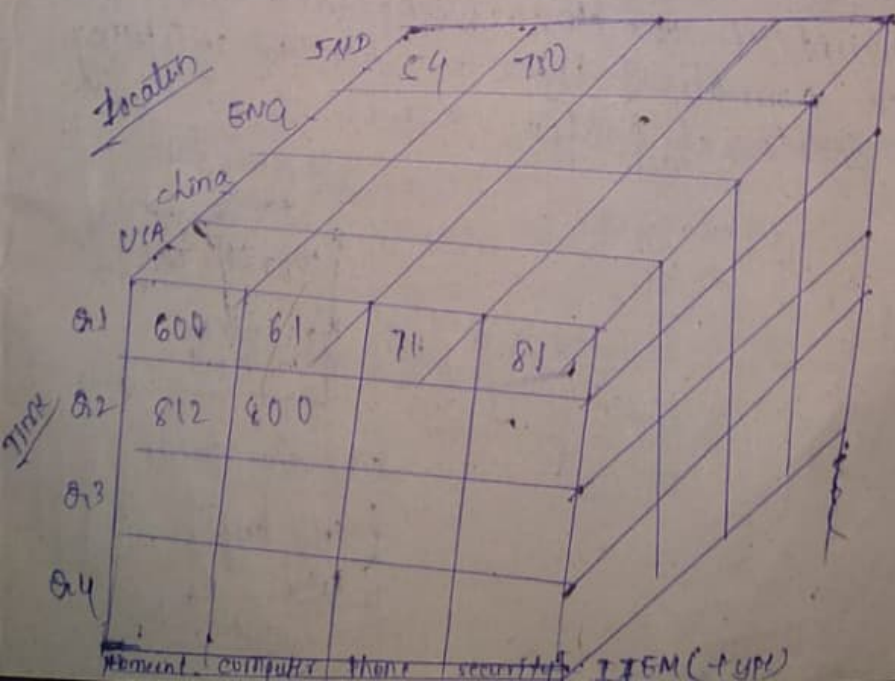
- Views data in form of data cube.
- Data cube - Data to be modelled and viewed in multiple dimensions.
- Defined by dimensions & facts.
- Dimensions → time, item, branch
- Facts → Dollars sold, units sold

2D view of sales data

Location USA

| item (type) |           |          |       |          |
|-------------|-----------|----------|-------|----------|
| time        | home ent. | Computer | Phone | Security |
| Q1          | 605       | 825      | 14    | 200      |
| Q2          |           |          |       |          |
| Q3          |           |          |       |          |
| Q4          |           |          |       |          |

3D view of sales data

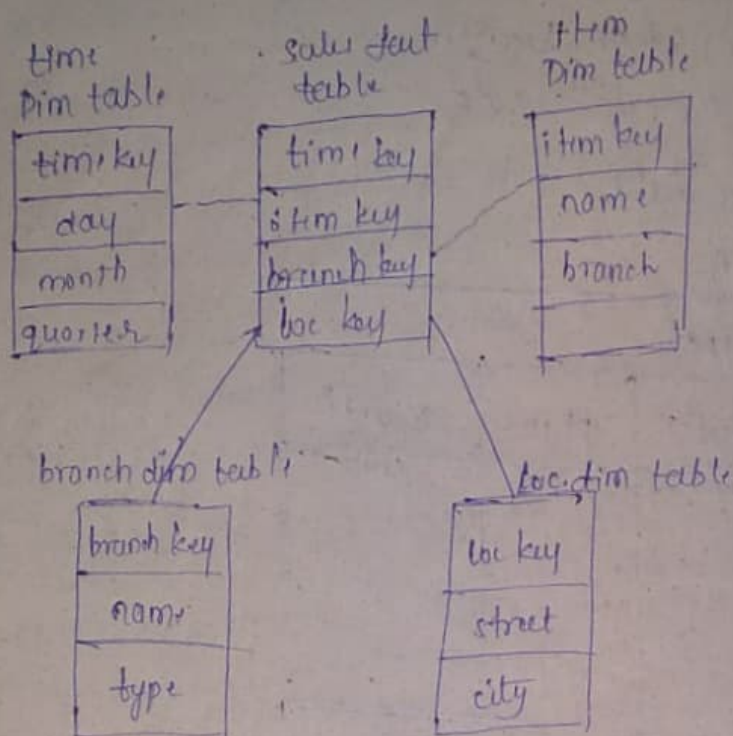


star, snowflake and fact constellations.

Schema for multidimensional DB (10-11 M 4)

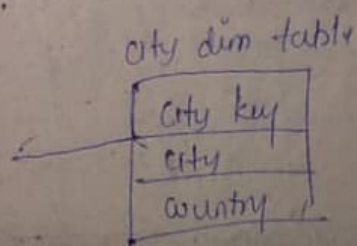
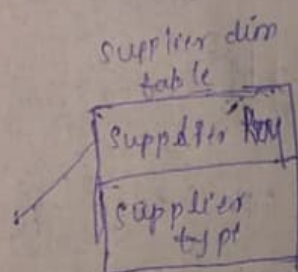
1) star schema of Datawarehouse for sales

Def<sup>2</sup>: The Data warehouse contains a large center table (fact table) containing the bulk of data, with no redundancy, set of smaller attribute table (dimension table) one for each dimension. (Many dim table one fact table)



2) Snowflake Schema

It is the variant of the star schema table where some dim tables are normalized there by further splitting the data into additional tables.





### 3) Fact constellation schema

It requires multiple fact tables to share dimension tables.  
This kind of schema can be viewed as collection of star schemas.  
It is called as galaxy schema or fact constellation schema.  
(Many fact table one dim table)

sales fact table

|          |
|----------|
| time key |
| item key |
|          |

shipping fact table

|              |
|--------------|
| item key     |
| time key     |
| shipping key |

Examples for defining star, snowflake and fact constellation schema

- SQL based ~~DM~~ query language - DML
- DW defined using 2 language primitives
- one for cube def<sup>2</sup>
- one for dimension def<sup>2</sup>

Syntax Define cube <cube name> <dimension list> : (<measure list>)  
Define dimension <dimension name> as (<attribute list>)

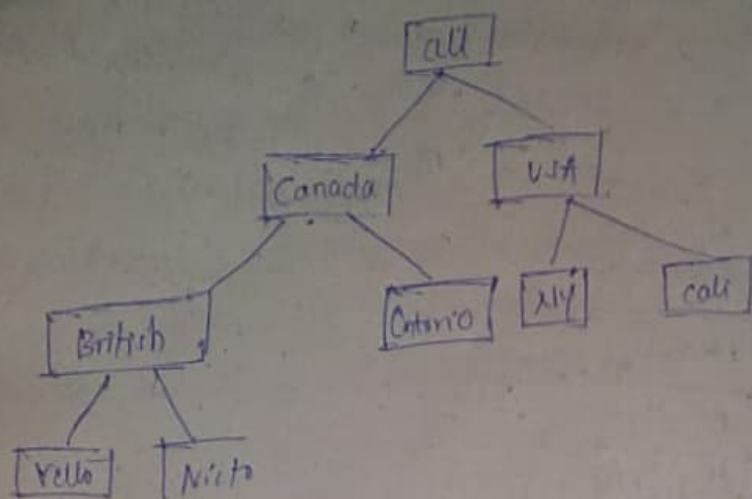
e.g., star schema.

define cube sales-star [time, item, branch, location] :  
dollars-sold = sum (sales-in-dollar), units-sold = count (\*)  
define dimension time as (time-key, day, month, year)  
" " item as (item-key, name, brand)  
" " branch as (branch-key, items)

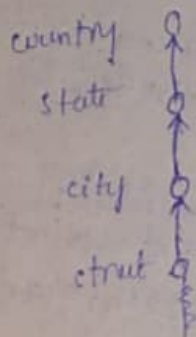
e.g., snowflake schema

define cube sales-star [time, item, branch, location, supply type]  
define dimension location as (location-key, street, city (city-key), country)  
define dimension item as (item-key, name, type, supplier (supplier type))

## Concept Hierarchy



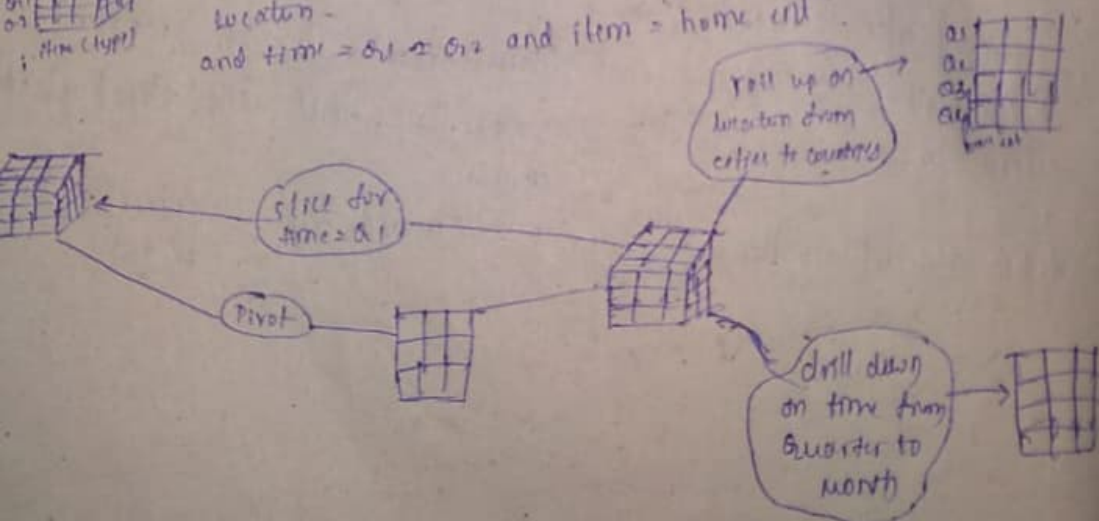
def 2: Concept hierarchy defines a sequence of mappings from a set of low level concepts to higher level. The mappings form a concept hierarchy for the dimension location mapping a set of low level concepts (cities) to higher level concepts (countries)



OLAP operation in multidimensional data model 07/12/2023



drill for location - and time =  $o_1 \pm o_2$  and item = 'home url'



Roll up operation: climbing up a concept hierarchy for a dimension.

e.g. street, city, state, country

Drill down: stepping down a concept hierarchy for a dimension.

e.g. month, week, day

slice operation: Performs selection on 1 D of the given cube

e.g. time = Q1

Pice operation: defines subcube by performing a selection on two or more dimensions.

Pivot: It provides an alternative presentation of data.

e.g. Rows  $\rightarrow$  columns or columns  $\rightarrow$  rows.

## Datawarehouse Architecture

$\rightarrow$  steps for design and construction of DW

What does DW provide for business analyst

- 1) Presents relevant information.
- 2) Enhance business productivity.
- 3) Customer relationship management
- 4) Cost Reduction

## A different views regarding the design of data warehouse

- 1) Top-down views: Allows the selection of relevant info necessary for the DW.
- 2) Data stores view: The info being capture, stored and managed by operational systems.
- 3) Data warehouse view: includes fact and dimension table
- 4) Business querying view: The viewpoint of the end user



## Database Design Process Steps

- 1) Choose a business process to <sup>be</sup> modeling, ordering, shipping, etc.
- 2) Choose grain of the ~~result~~ business process key, Individual transactions
- 3) Choose the dimensions.
- 4) Choose the measures for each fact table record.

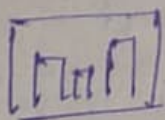
20/12/2023

## Process of Database Design

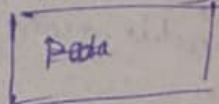
- Top down approach - design of planning.
- Bottom-up approach - experiment and prototypes.
- Combined approach.
- Software engg → Planning requirement study, problem analysis, warehouse design, Data integration and testing deployment.

## Three-Tier DW Architecture (8-10 M) \*\*\*

Query / Report



Data mining



Top-Tier  
front end  
tools.

OLAP server

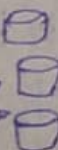


Middle Tier  
OLAP server

Administration

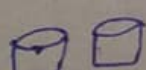
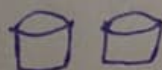


DW



Data mart

Bottom Tier  
DW server



External  
sources

\* Bottom tier is almost a relational DB where the data extracted from the tier in order to create a data warehouse.

\* Middle tier is implemented using relational OLAP and multidimensional OLAP.

\* Top tier contains reporting tools, analysis tools and data mining tools.

### Enterprise <sup>Data</sup> Warehouse and Data Mart

→ Collects all of the information about subjects spanning the entire organization.

Data Mart: It contains information specific to an organization business unit.

### Types of OLAP servers

- 1) Relational OLAP server.
- 2) Multidimensional OLAP server.
- 3) Hybrid OLAP server.

#### 1) Relational OLAP server

→ Application based on relational DBMS. It can handle large amount of info., it has greater scalability, and its performance is slower low.

#### 2) Multidimensional OLAP server

Application based on multidimensional DBMS. It has fast info. retrieval, performance complex calculation and limited info. it can handle.

→ Through indexing info. can be retrieved.

#### 3) Hybrid OLAP server

→ It is combination of both relational and multidimensional OLAP server.



Datwarehouse ImplementationCompute cube operator and its implementation

→ Compute cube → Aggregate over all subsets of dimensions.  
e.g., 3 dimensions - city, item, year and sales in dollars as measures

$2^3 = 8$  Possible groupings {city, item, year}, {city, item}, {city, year}, {item, year}, {city}, {item}, {year}, {}

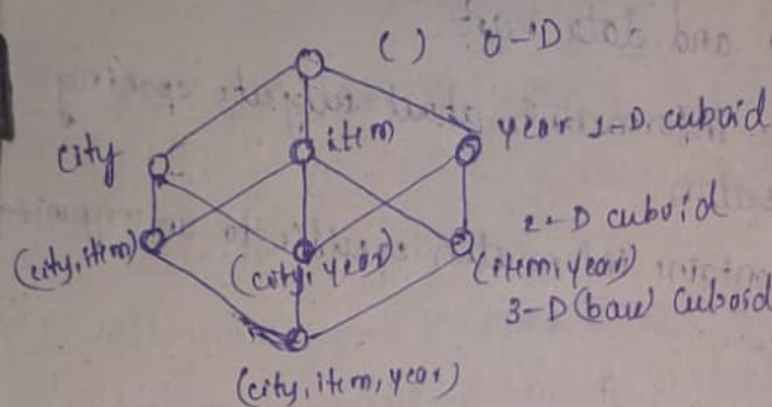


Fig. 3-D Database.

Datwarehouse Backend Tools and Utilities.

1. Data extraction
2. Data cleaning
3. Data Transformation
4. Load
5. Refresh → All the updates to the data has to be done

## UNIT-04

### CLUSTERING

DEF \*\*\*  
(2-10 M)

Cluster is a collection of data objects that are similar to one another within the same cluster and dissimilar in other cluster.

→ It is widely used in numerous applications like image processing, pattern recognition, data analysis, market research.

→ Clustering is an example of unsupervised learning.

\* Unsupervised learning do not rely on predefined classes and class labeled training examples.

\* Supervised → Predefined sets or classes.

### Requirements of clustering

1. Scalability
2. Ability to deal with different types of attributes.
3. Minimal requirements for domain knowledge to determine input parameters.
4. Ability to deal with noisy data.
5. High dimensionality.
6. Constraints based clustering.
7. Interpretability and stability.

### Types of data in cluster analysis

1. Interval scaled variable → 200-300
2. Binary variable → 0 or 1
3. Nominal variable → map, colour (representing more than one state)
4. Ordinal → Assistant, Associate, professor (order/sequence)



## Categories of clustering methods

### 1) Partitioning method

- DB of  $n$  objects,  $k$  partitions of data
- $k \leq n$
- classify data into  $k$  groups
  - 1) Each group must contain at least one object
  - 2) Each object must belong to exactly one group
- Uses iterative relocation technique methods
  - 1)  $k$ -means
  - 2)  $k$ -medoid

### 2) Hierarchical Method

- Hierarchical decomposition of given set of data object
- Methods
  - 1) Agglomerative → Bottom-up
  - 2) Divisive → Top-down

### 3) Density Based Method

- It is based on density, it continuously growing the given cluster as long as the density exceeds some threshold.

### A) Grid Based Method

- It quantizes the object space into finite no. of cells that form a grid structure. All the clustering operations are performed on cluster grid structure.



## Partitioning Method (K-Means)

→  $K \leq n$

→ K-means Algorithm

I/P: no. of clusters  $k$ ,  $n$  objects.

O/P: set of  $k$  clusters.

### Method

- 1) Arbitrarily choose  $k$  objects as initial cluster centers.
- 2) Repeat
- 3) (a) assign each object to cluster to which the object is most similar based on mean value of the object in the cluster.
- 4) Update the cluster mean.
- 5) Until no change.

$k = 3$



### How k-means algorithm works:

- First it randomly selects  $k$  of the objects, each of which initially represents cluster means or cluster centre.
- For each of the remaining objects an object is assigned to a cluster to which it is most similar based on the distance b/w object and cluster mean.

then compute the new mean for each cluster. This process iterates until all the objects are in one of the  $k$  clusters.

### K-medoid Algorithm

Input: No. of  $k$  clusters,  $n$  objects

Output: Set of  $k$  clusters.

#### Method:

- 1) Arbitrarily choose  $k$  objects as initial medoids.
- 2) Repeat
- 3) Assign each remaining object to cluster with nearest medoid.
- 4) Randomly select non medoid object  $O_{rand}$
- 5) Compute total cost,  $c$  of swapping  $O_j$  with  $O_{rand}$
- 6) If  $c < 0$  then swap to form new set of  $k$  medoids
- 7) Until no change.

### K-medoid Algorithm : (8M) \*\*\*

|                                   | X | Y | k <sub>1</sub> | k <sub>2</sub> | cost |
|-----------------------------------|---|---|----------------|----------------|------|
| k <sub>1</sub> = m <sub>1</sub> ① | 8 | 2 | 0              | 6              | 0    |
| ②                                 | 3 | 5 | 8              | 2              | 9    |
| ③                                 | 4 | 7 | 9              | 3              | 3    |
| ④                                 | 8 | 4 | 2              | 4              | 2    |
| k <sub>2</sub> = m <sub>2</sub> ⑤ | 5 | 5 | 6              | 0              | 0    |
|                                   |   |   |                | 7              |      |

#### Manhattan Distance

$$\rightarrow |x_2 - x_1| + |y_2 - y_1|$$

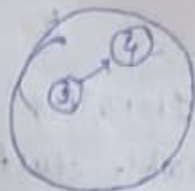
$$= |15 - 8| + |15 - 2| = 3 + 3 = 6$$

$$\rightarrow |13 - 8| + |15 - 2| = 5 + 3 = 8$$

$$\rightarrow |13 - 5| + |15 - 5| = 2$$



$$\begin{aligned} \rightarrow |4-8| + |7-2| &= 4+5=9 \\ \rightarrow |4-5| + |7-5| &= 1+2=3 \\ \rightarrow |8-8| + |4-2| &= 2 \\ \rightarrow |8-5| + |4-5| &= 3+1=4 \end{aligned}$$

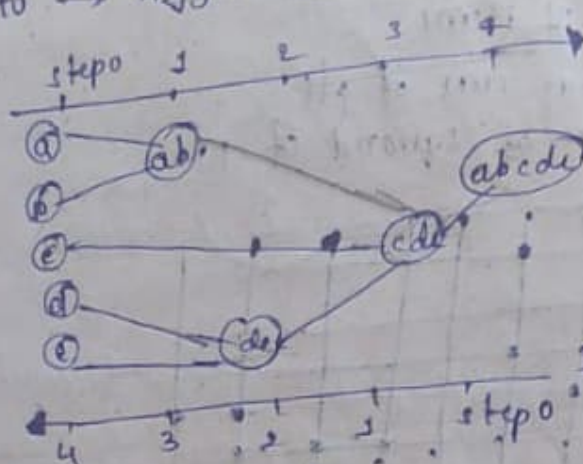


04/01/2024

## Hierarchical Method (\*\*\* try)

- Works by grouping data objects into tree of clusters.
- Divided into → Agglomerative & Divisive Hierarchical clustering.

Agglomerative



Divisive

## Agglomerative Hierarchical clustering:

This is bottom-up strategy starts by placing each object in its own cluster and then merge these atomic clusters into larger & larger clusters until all the objects are in single cluster.

## Divisive Hierarchical clustering:

This is top-down approach it starts with all objects in one cluster. It subdivides the cluster into smaller & smaller pieces until each object forms a cluster on its own.



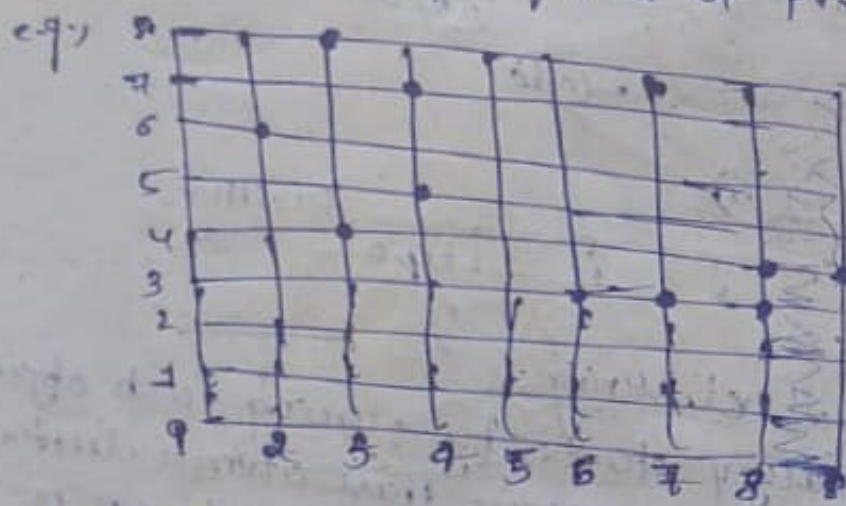
## Balanced Iterative Reducing & Clustering using Hierarchical

- It is scalable clustering.
- Works for very large dataset.
- Only one scan of data is necessary.
- Clustering is based on CF (Clustering Feature).
- CF Tree → stores the cluster features.
- cluster of data pts is represented by triple of numbers  $(N, L, SS)$

$N$  = No. of items

$L$  = Linear sum of pts.

$SS$  = Sum of squares of pts.



e.g.) CF

$$CF = (N, L, SS)$$

$N$  = No. of data pts

$$L = \sum_{i=1}^N X_i \quad SS = \sum_{i=1}^N X_i^2$$

$$(3, 4), (2, 6), (4, 5), (4, 7), (3, 8)$$

$$N = 5$$

$$L = 3 + 2 + 4 + 4 + 3 = 16$$

$$= 4 + 6 + 5 + 7 + 8 = 30$$

$$SS = 3^2 + 2^2 + 4^2 + 4^2 + 3^2 = 54$$

$$4^2 + 6^2 + 5^2 + 7^2 + 8^2 = 190$$

(6, 7) (7, 8) (8, 3) (8, 4)

3)  $N = 4$

$$CF = (M, LS, SS)$$

$$LS = 6 + 7 + 1 + 8 = 29$$

$$2 + 3 + 3 + 4 = 12$$

$$SS = 6^2 + 7^2 + 1^2 + 8^2 = 213$$

$$2^2 + 3^2 + 3^2 + 4^2 = 36$$

$$CF = (4, \dots)$$

$$\begin{array}{r} 30 \\ 49 \\ 64 \\ \hline 143 \end{array}$$

Basic algorithm of BIRCH: Phase 1

Phase 1: Load the data into memory.

Phase 2: Condense data - resize the data set by building smaller cluster feature (CF) tree.

Phase 3: Global clustering - It uses existing clustering algorithm on CF entries.

Phase 4: Cluster refining.

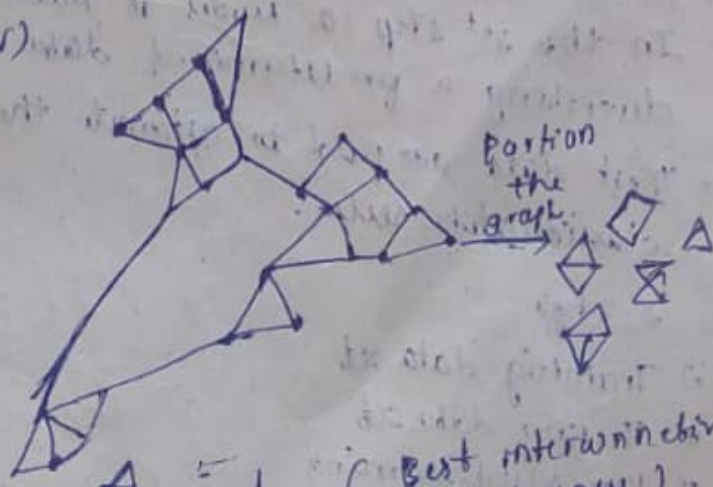
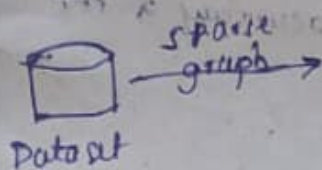
Chameleon: Hierarchical clustering using dynamic modeling

→ Two clusters are merged → if interconnectivity & closeness b/w 2 clusters are high.

→ 2 parameters → Interconnectivity & closeness.

→ Graph based & a two phase algorithm graph partitioning algorithm.

→ (ii) is agglomerative hierarchical clustering algo. Construct (K-NN)



Final cluster (Best interconnectivity and closeness)

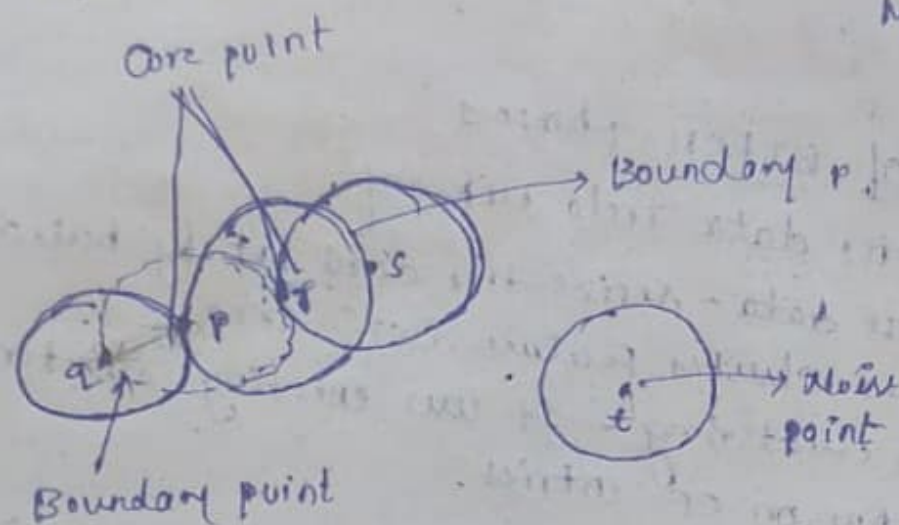


(Density Based Spatial clustering of Application with Noise).

→ Forms clustering based on density.

→  $\epsilon$  p1  
Minpts = 3 ] → 2 t/p parameters

Minpts → Minimum pts.



→ Directly density reachable must be neighbour of core pt.  
so q is directly density reachable from pt p.

→ p, q, r are core pts → satisfies Minpts = 3

→ s, q, r → Boundary pt → should be neighbour of any core point. t is noise point.

## Classification and Prediction

→ Data classification is a 2 step process.  
In the 1st step a model is built describing a predetermined data class or concepts.  
→ Test data are used to estimate the accuracy of the classification rules.

### 2 steps

- Training data set
- Testing data set
- Supervised learning
- Unsupervised learning





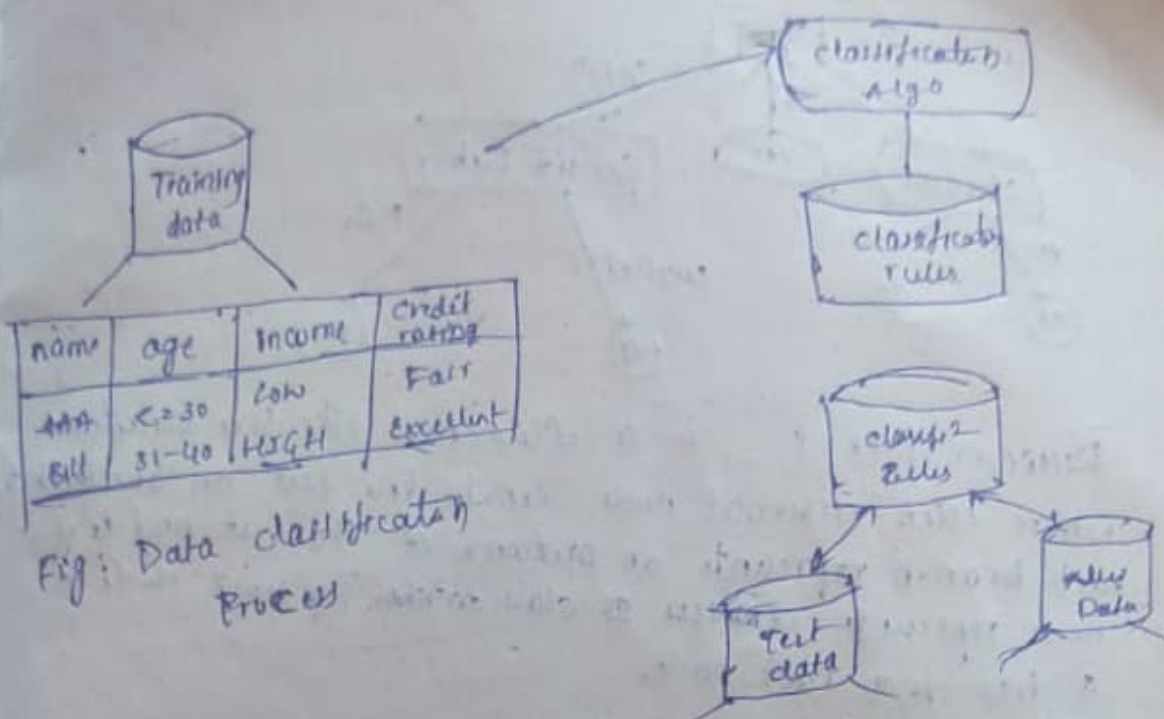


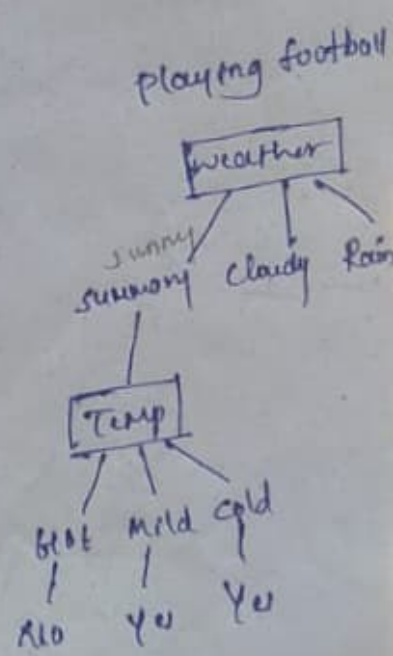
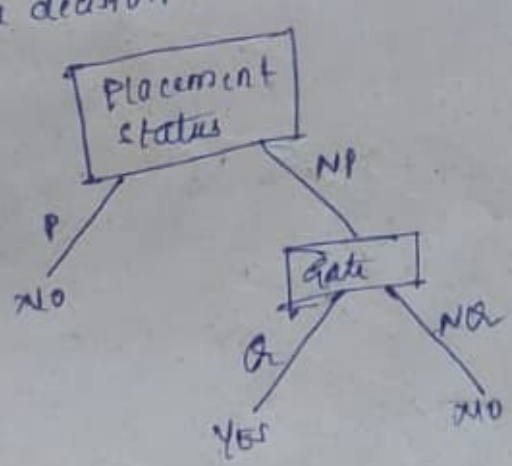
Fig: Data classification Process

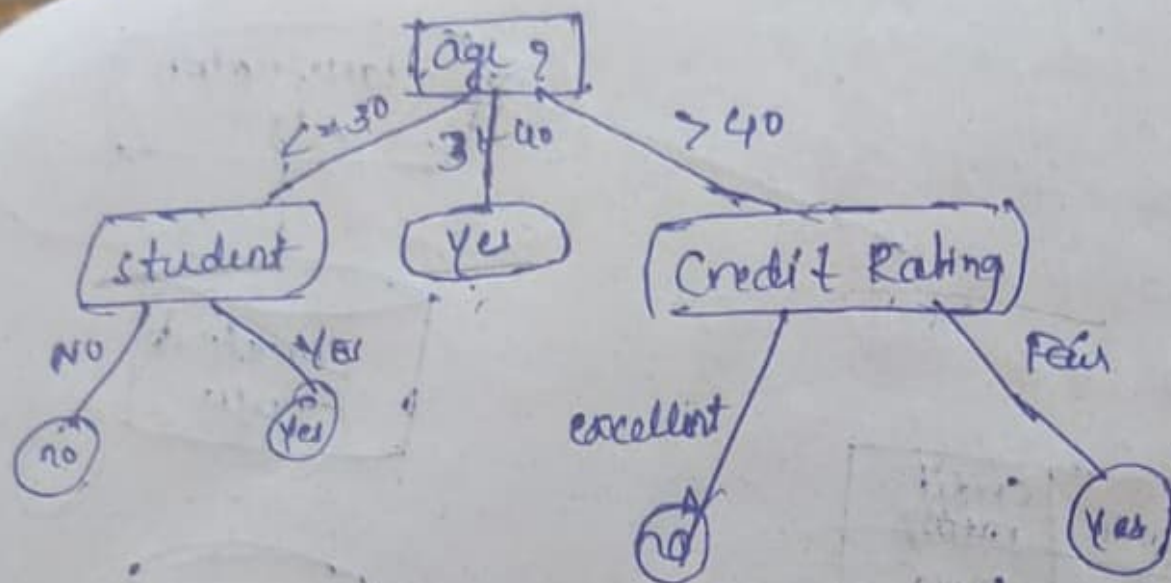
## Issues regarding classification and prediction.

- 1) Data cleaning: Preprocessing of data to remove or reduce noise.
- 2) Relevance analysis.
- 3) Data Transformation: Data can be generalized to higher level concepts.

## Decision Tree Induction

→ To take a decision





**Decision Tree :** is a flowchart like tree structure where each internal node denotes the test on an attribute, each branch represents an outcome of the test and leaf node represents classes or class distribution. Top most node in a tree is a root node.