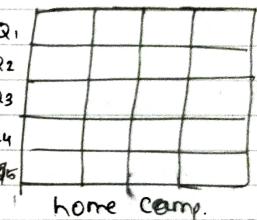
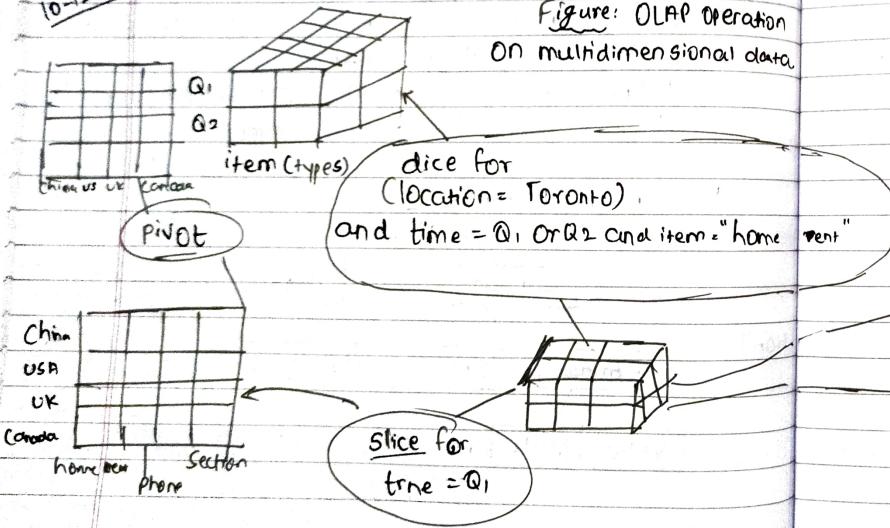


- 1) Roll up Operation: Climbing up a concept hierarchy @ for a dimension  
Eg: Street , City, State Country.
- 2) Drill down Operation: Stepping down a concept hierarchy for a dimension.  
Eg: Time

### \* OLAP Operations in multidimensional data model

10-12 marks



5) Slice and dice: Slice operation performs a selection on one dimension of given cube.

4) Dice operation defines a sub cube by performing a selection of 2 or more dimensions.

5) Pivot: It provides an alternative presentation of data.

### \* Datawarehouse Architecture

→ Steps for a design Of construction Of DW

- What does DW provide for business analyst?
- i) Present relevant information?
- ii) Enhance business productivity
- iii) Customer relationship management
- iv) Cost Reduction.

The 4 different views regarding the design of datawarehouse

1) Top-down View: It allows the selection of relevant info necessary for datawarehouse

2) ~~data warehouse~~ <sup>source</sup> view: The info captured, stored and managed by systems Operational Systems.

3) Datawarehouse view: Includes fact tables & dimensional table

4) Business query: This is the view of the end user.

### \* Datawarehouse Design Process Steps

1) Choose a business process to model  
Eg: Orders, Shipments, Sales.

2) Choose grain of the business process - Eg: Individual transaction

3) Choose the dimensions

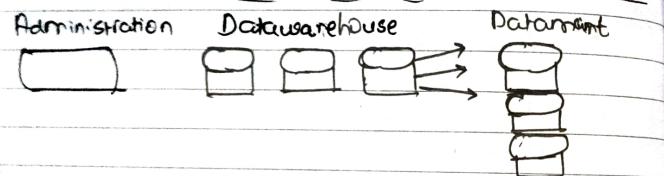
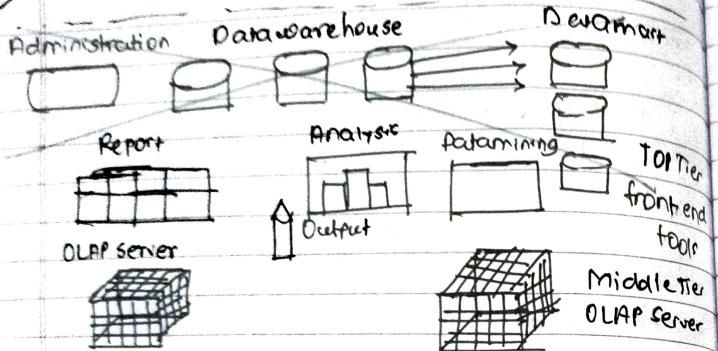
4) Choose the measures for each fact table record  
Eg: Numerical measures such dollars exchanged

### Process Of Datawarehouse Design

- TOP down approach → Design and planning
- Bottom up approach → Experiment And prototypes
- Combined approach

Software Engineering → planning requirement study, problem analysis, warehouse design, Data integration and testing, deployment.

## Three-Tier Datawarehouse Architecture



- 1) Bottom Tier is almost a relational database where the data extracted from this tier in order create a data warehouse.

2) Intermediate Middle Tier: is implemented using relational OLAP and multidimensional OLAP.

3) Top Tier: Contains reporting tools, analysis tools and DM tools

4) Enterprise Warehouse and Datasmart: It collects all of the information about subjects spanning entire organization.

Datasmart: It contains information specific organization business unit.

### Types Of OLAP Servers

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

1) Relational OLAP server: Applications based on relational DBMS. It can handle large amount of info, greater scalability, performance is slow.

2) Multidimensional OLAP server: Applications based on multidimensional DBMS, it has fast info retrieval, complex calculations is performed and limited info it can handle.

### 3) Hybrid OLAP Server:

#### Datawarehouse Implementation

Compute cube Operator And it's implementation  
→ Compute cube → Aggregates Overall Subsets Of dimensions.

Eg:

3 dimensional - City, item, year And sales in dollar as measure

$2^3 = 8$  possible groupby - {city, item, year}  
{city, item}, {city, year}, {item, year},  
{city}, {item}, {year}

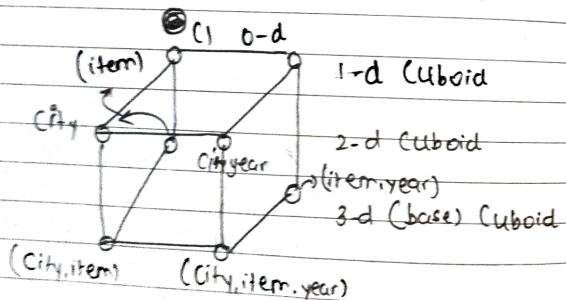


Figure: 3d datacube

#### \* Datawarehouse Back-end tools & Utilities

- 1) Data Extraction
- 2) Data Cleaning
- 3) Data Transformation
- 4) Load
- 5) Refresh.

21-12-23

## Unit-4 Clustering

Viva  
80-100%

\* Cluster: It is a collection of Objects which is similar to one another and dissimilar to objects in other clusters. It is widely used in pattern recognition, image processing, and market research.

Clustering is an example of UnSupervised learning.

Unsupervised learning don't depend on predefined classes and class labelled training example.

### 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) Constraint based clustering
- 7) Interpretability and Usability.

### Types Of Data in Cluster Analysis

- 1) Internal Scaled Variable  $\rightarrow$  200-300
- 2) Binary Variable  $\rightarrow$  0 or 1
- 3) Nominal Variable  $\rightarrow$  Mar. Status, color
- 4) Ordinal Variable  $\rightarrow$  Assistant, Associate Professor

Categories of clustering methods

#### 1) Partitioning method

$\rightarrow$  DB Of N Of Objects, K portions Of data  
 $\rightarrow$   $K \in \mathbb{N}$

$\rightarrow$  classification classify data into k groups  
1) Each group must contains atleast 1 object  
2) Each Object must belong to each group.

$\rightarrow$  Uses iterative relocation technique

$\rightarrow$  Methods

- 1) K-means
- 2) K-medoid

#### 2) Hierarchical Method

$\rightarrow$  Hierarchical decomposition Of given set of data

Objects

$\rightarrow$  Methods

- 1) Agglomerative  $\rightarrow$  Bottom up approach
- 2) Divisive  $\rightarrow$  Top down approach

\* Density based method: It is based on density, it continues growing given cluster as long as the density exceeds some threshold.

\* Grid based method: It quantizes the object space into finite space number of cells that form a grid structure. All clustering operations performed on grid structure.

\* Partitioned based Method:

$$\rightarrow K \leq n$$

\* K-means algorithm (Centroid based Techniques)

I/P: No of clusters  $K$ ,  $n$  objects

O/P: Set of  $K$  clusters

Method:

1) Arbitrarily choose  $K$  objects as initial cluster centres

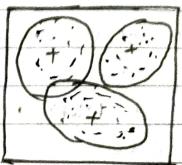
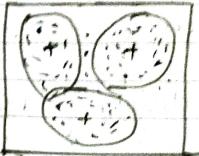
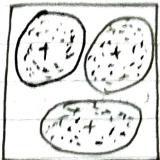
2) Repeat

3) (re) assign each of object to cluster to which the object is most similar based on the value of objects in the cluster.

4) Update the cluster

5) Until no change

$$K=3$$



\* How K-means algorithm works?

① First it selects randomly select  $K$  of the objects each of which represents cluster mean and cluster centre. For each of the remaining objects assign ~~an object~~ to the cluster to which it is the most similar based on the distance between the object and the cluster mean. It then computes new mean for each cluster. This process iterates until all the objects are inside a cluster.

## \* K-medoid algorithm

I/P Number Of K Clusters , n Objects

O/P Set Of K clusters

Method:

- 1) Arbitrarily choose K objects as initial medoid
- 2) Repeat
- 3) Assign each remaining object to cluster with nearest medoid
- 4) Randomly select non-medoid objects or random
- 5) Compute total cost. S of swapping  $O_j$  with random
- 6) If  $S \leq 0$  then swap to form new set of K-medoids
- 7) Until no change.

## K-medoid algorithm

	X	Y	$K_1$	$K_2$	COST
①	5	2	0	6	80
②	3	5	8	2	2
③	4	7	9	3	3
④	8	4	2	4	2
⑤	5	6	0	0	0
					7

$K_1 = m_1$

$K_2 = m_2$

## Manhattan Distance

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

$$\rightarrow |5-8| + |5-2| = 3+3 = 6$$

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

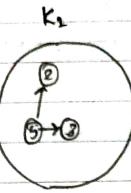
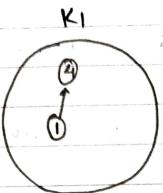
$$\rightarrow |3-5| + |5-5| = 2$$

$$\rightarrow |4-8| + |7-2| = 4+5 = 9$$

$$\rightarrow |4-5| + |7-5| = 1+2 = 3$$

$$\rightarrow |8-8| + |4-2| = 0+2 = 2$$

$$\rightarrow |8-5| + |1-0| = 3+1 = 4$$

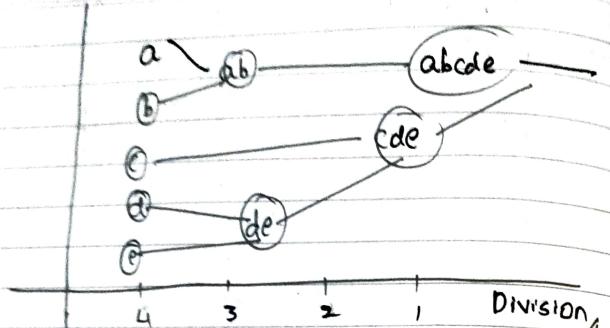


## Hierarchical Method

→ Works by grouping data objects into tree of clusters

→ Dividend into → Agglomerative and divisive hierarchical clustering. clusters.

## Agglomerative



- 1) This is top down approach it starts with all objects in one cluster.
- 2) Starts with all objects & it sub divides cluster into smaller and smaller pieces until each objects form cluster on its own.

\* Balanced Iterator reducing and clustering using Hierarchies

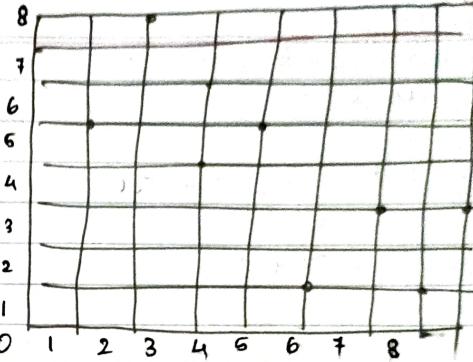
- It's scalable clustering
  - Works for very large data sets
  - Only one scan of data is necessary feature
  - Clustering is based on CF (Cluster Feature)
  - CF Tree → Stores the Cluster feature
- Cluster of data pts is represented by triple of numbers  $(N, LS, SS)$

$N = \text{No of Iterations}$

$LS = \text{Linear sum of points}$

$SS = \text{Sum of squared of pts}$

Eg:



Eg: CF

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

$N = \text{No of data pts}$

$$LS = \sum_{i=1}^n X_i$$

$$SS = \sum_{i=1}^n X_i^2$$

$$(2, 6), (3, 8), (4, 4), (4, 6), (5, 6), (6, 1), (6, 6), (7, 3), (8, 2), (8, 4)$$

$$(8, 4)$$

$$N = \sum_{i=1}^n X_i = 2 + 3 + 4 + 4 + 5 + 6 = 29$$

$$SS = \sum_{i=1}^n X_i^2 = 2^2 + 3^2 + 4^2 + 4^2 + 5^2 + 6^2 + 8^2$$

$$= 4 + 9 + 16 + 16 + 25 + 36 + 64$$

$$= 170.$$

LS:  $\sum_i^N y_i = 6+8+4+6+6+0+4$   
 $= 34$

SS =  $\sum_i^N y_i^2 = 6^2 + 8^2 + 4^2 + 6^2 + 6^2 + 0^2 + 4^2$   
 $= 36 + 64 + 16 + 36 + 36 + 16$   
 $= 204.$

### Basic Algorithm Of BIRCH

#### Phase①

① Load the data into memory

#### Phase②:

Condense Data - Resize the data set by building smaller cluster feature tree

#### Phase③:

Global Clustering- It uses very existing clustering algorithms on CF entries

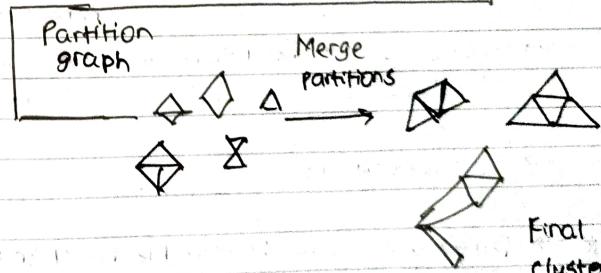
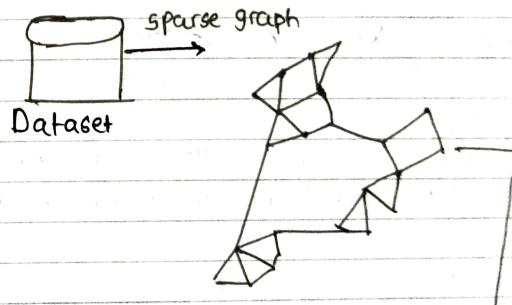
#### Phase④:

Cluster Refining

\* Chameleon: Hierarchical Clustering using dynamic modelling.

- Two clusters merged → if interconnecting or closeness b/w 2 clusters are high
- 2 parameters → Interconnectivity And Closeness.
- Graph based and a two phase algorithm.
  - Uses a graph partitioning algorithm
  - Uses agglomerative hierarchical cluster features.

#### Construct (K-NN)

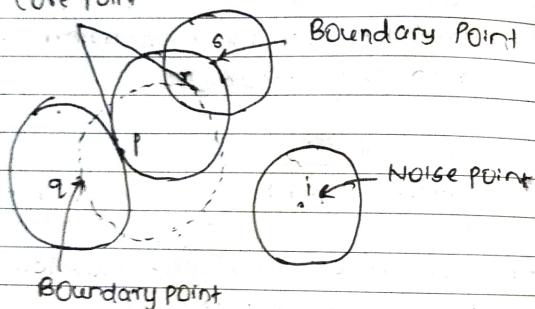


## ★ DBSCAN

Vineet  
+ friends  
(Density based Spatial Clustering Of application with noise)

- Forms Clustering based On density
- EPS → 2 i/p Parameters
- Mmpts = 3

Core Point



→ P and Q Core Points → satisfies  
 $m m p t s = 3$

→ S and Q → Boundary points → Should  
be neighbour neighbour Of any core point.

→ T is noise pt

★ Directly density Reachable must be neighbour  
and Core point SO Q is directly density reachable  
from pt P.

5/11/24

## Unit-3 Classification and Prediction

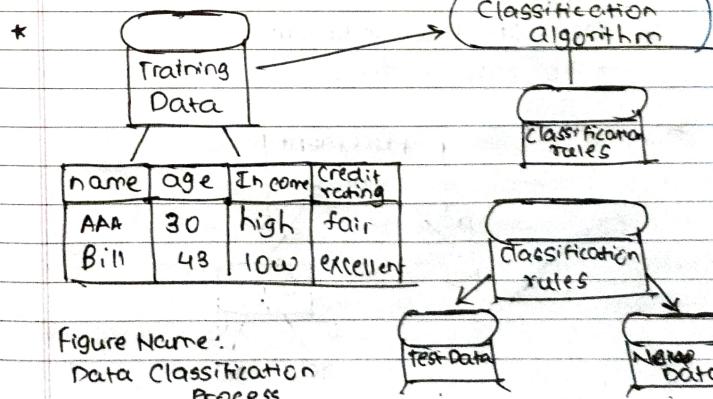
Classification: Data classification is 2 step process

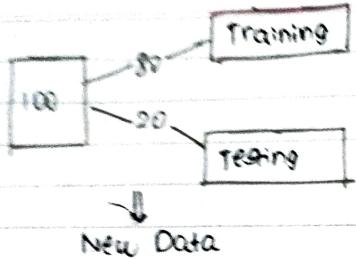
- 1) A model is built determining describing a pre determined data classes or concepts

- 2) Test data are used to estimate accuracy of the classification rules.

## Classification

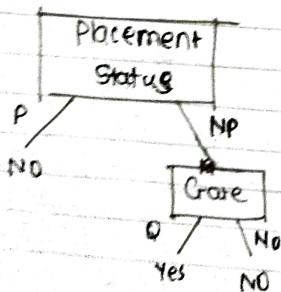
- Training dataset
- Testing dataset
- Supervised Learning → labelled data
- Unsupervised Learning → unlabelled data



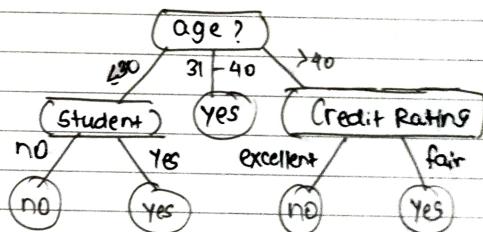
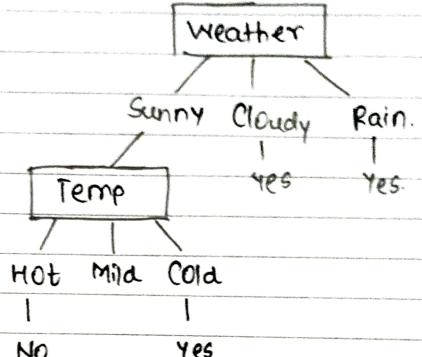


Issues regarding Classification and Prediction

- 1) Data Cleaning: Pre processing Of data to reduce noise
- 2) Relevance Analysis: Reducing irrelevant attributes
- 3) Data Transformation: Data can be generalized to higher level concept
- 4) Decision Tree Induction  
→ To take a decision.



Playing Football



- \* Decision tree is a flowchart like tree structure where each internal node denotes a test on an attribute , each branch denotes represents the outcome of test and leaf node represents class or class distribution. TOP most node of a tree is the root node .