

**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.

0410112024.

Abair)

—2 Ep 7

2/18/2008

maps  $\rightarrow$  Minimax  
pts.



longer in 155

$\rightarrow$   $P \leq r$  or  $r \leq P$

any core point, it is not a core

## Classification and Infection

→ Data classification is a & set of built described of

describing a phenomenon and the accuracy of the

Cavitation rules:

2559

→ Training data set

2) Testing data set

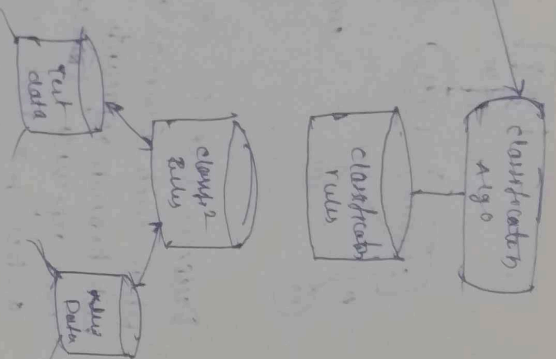
→ Super Wild Morning

→ unsupervised learning

name	age	income	credit rating
A1A4	<=30	low	Fair
B11	31-40	High	Excellent

Fig: Data cross-section

Process



Issue regarding classification -

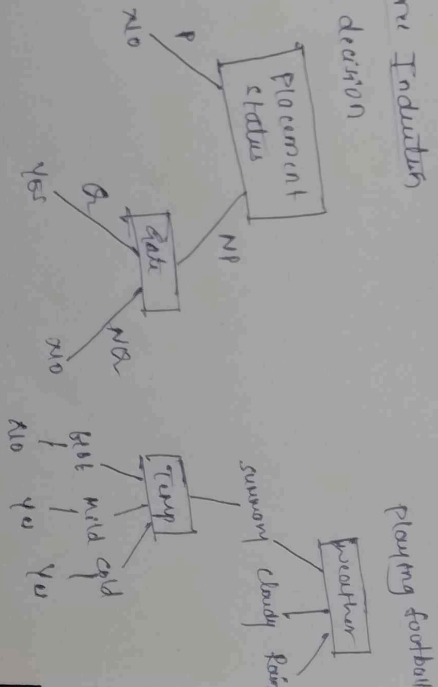
17 Data cleaning: info

→ Relevant analysis

concepts.

## Diagon In Inductor

→ To take a decision



$\epsilon$

$1100$

$(1) \rightarrow 100$

$\text{long } k+1$

$m=1 \rightarrow M$

$n=1$

$11$

$11$

$001$

$[1]$



## BIRCH

(Balanced Iterative Reducing & Clustering using thresholds)

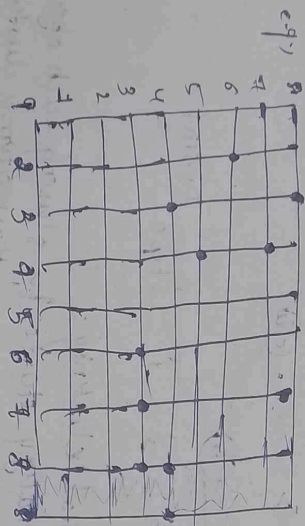
- It is scalable clustering.
- Works for very large dataset.
- Only one scan of data is necessary.
- Clustering is based on CF (Clustering Features).
- CF Tree → store the cluster features.
- cluster of data pts is approximated by triple of numbers

$(C_i, L_i, SS)$

$N_i$  = No. of items

$L_i$  = Linear sum of pts

$SS$  = Sum of squared of pts.



e.g)

CF =  $(C_i, L_i, SS)$

$N_i$  = No. of data pts

$L_i = \sum_{j=1}^N X_j$   $SS = \sum_{j=1}^N X_j^2$

$(3, 41) (4, 6) (4, 5) (4, 7) (3, 8)$

$N = 5$

$L = 3 + 2 + 4 + 4 + 8 = 21$

$SS = 3^2 + 2^2 + 4^2 + 4^2 + 8^2 = 89$

$SS = 3^2 + 2^2 + 4^2 + 4^2 + 8^2 = 89$

$SS = 3^2 + 2^2 + 4^2 + 4^2 + 8^2 = 89$

e.g)  $(4, 3) (4, 3) (4, 4)$

$N = 4$

CF =  $(N, L, SS)$

$L = 4 + 7 + 1 + 8 = 29$

$SS = 4^2 + 7^2 + 1^2 + 8^2 = 122$

$SS = 4^2 + 7^2 + 1^2 + 8^2 = 122$

CF =  $(4, 29, 122)$

Basic algorithm of BIRCH:

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: Clustering refining.

Character: hierarchical clustering using dynamic modeling.

Two clusters are merged → if interconnectivity & closeness.

→ 2 partitions → Interconnectivity & closeness.

→ Graph based & a two phase algorithm.

→ Graph algorithm.

→ (iii) Use aggregative hierarchical clustering algo.

Construct (K=N) graph.

Dataset

space graph

partition

Find clusters (cut & interconnectivity)

Find clusters (cut & interconnectivity)

Find clusters (cut & interconnectivity)

Find clusters (cut & interconnectivity)

Find clusters (cut & interconnectivity)

Find clusters (cut & interconnectivity)

Find clusters (cut & interconnectivity)

Find clusters (cut & interconnectivity)

→ then compute the new medoid for each cluster. This process iterates until all the objects are in one of the clusters.

### K-medoid Algorithm

ILP: No. of k clusters, n objects.  
OP: 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  $C_{random}$
- 5) Compute total cost,  $c$  of swapping  $C_j$  with  $C_{random}$
- 6) If  $c < 0$  then swap to form new set of k medoids
- 7) Until no change.

#### k-medoid Algorithm

(8N) \*\*\*

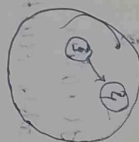
	X	Y	k1	k2	cost
k1 = m1	0	2	0	6	0
	3	5	8	1	9
	4	7	9	3	3
	8	4	2	4	2
	9	5	6	0	0
	5	5	7	7	7

k2 = m2

#### Manhattan distance

- $|x_1 - x_2| + |y_1 - y_2|$
- =  $|15 - 21| + |5 - 3| = 8 + 2 = 10$
- $|13 - 8| + |15 - 21| = 5 + 6 = 11$
- $|13 - 5| + |15 - 5| = 8 + 10 = 18$

- $|14 - 8| + |17 - 21| = 6 + 4 = 10$
- $|14 - 5| + |17 - 5| = 9 + 12 = 21$
- $|12 - 8| + |14 - 21| = 6 + 7 = 13$
- $|18 - 5| + |14 - 5| = 13 + 9 = 22$

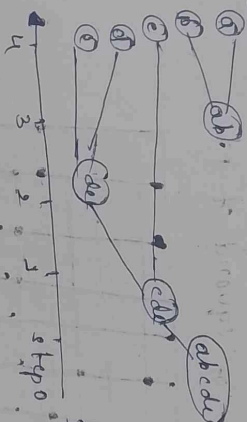


04/03/2023

### Hierarchical Method

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

#### Agglomerative

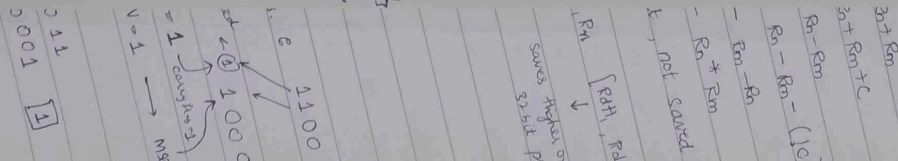


#### Agglomerative Hierarchical clustering

This is bottom-up strategy starts by placing each object in its own cluster and then merge the 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.





## Categories of clustering methods

## Pathoning Method

→ DB of  $n$  objects,  $k$  partitions of data

↓  
15  
5

→ classify data into  $k$  groups

1) Each group must contain atleast one object

a) Each object must belong to exactly one group

→ Uses iterative relocation technique methods

1)  $k$ -means

2) k-medium

## 2) Hierarchical Method

a) Hierarchical Method  
decomposition of given set of data objects

## Methods

1) Agglomeration  $\rightarrow$  Bottom-up

2) Division  $\rightarrow$  Top-down

3) Density Band Method

3) Density Base medium → It is based on density, it continuous growing the the density exceeds some

It is based on the density of cells in the given cluster as long as the density exceeds 50000.

Threshold.

A) Bald King Bard Method

A) Grid Based Method :-  
→ It quantifies 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 (10 SA) \* \* \*

六八五

→ K-means Algorithm

$k$  - no. of clusters  $k$ ,  $n$  objects.

$\frac{1}{O(p)} \cdot$  set of  $k$  clusters.

Method

Method  
1) Arbitrarily choose  $k$  objects as initial clusters

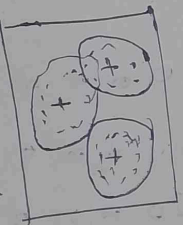
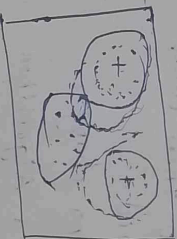
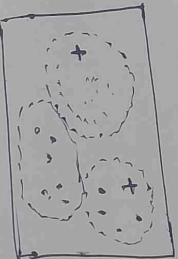
2) Report

- 2) Repeat
- 3) Core assign each object to cluster to which it is most similar based on mean

value of the word

- 4) Update the counter
- 5) Until no change.

75  
11  
6



k-means algorithm works

How k-means algorithm works:

- First it randomly selects  $k$  of the objects as cluster which initially represents cluster means or cluster centre.
- The remaining objects are assigned to similar based on

For each of the remaining  $n - 1$  clusters to a cluster to which it is most similar. the distance b/w centroid and cluster mean.



## Databarhouse Implementation

21/12/2023

Comput cube operator and its implementation

→ Comput cube → Aggregate over all subset of dimensions  
e.g. 3 dimensions - city, item, year and state-in-dollars as measures

$2^3 = 8$  Possible grouping city, item, year, city, item, year, year, item, year, city, item, year, year, year, year, year, year

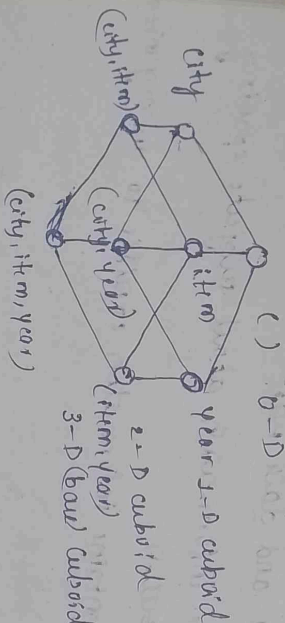


Fig. 3 D Databar.

## Databarhouse 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

DEFINITION

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 does not use on pre-defined 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 quality.

\* \* \*

### Type of data in cluster analysis

1. Interval scaled variable → 200-300
2. Binary variable → 0 or 1 (representing yes/no or state)
3. Nominal variable → map, color, profession (cardinal/several)
4. Ordinal → Assistant, Associate, Professor (cardinal/several)

## Database Design Process Steps

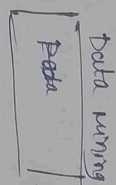
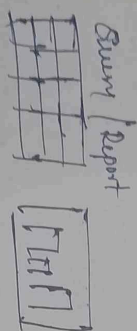
- 1) Choose a business process to model (e.g., order, shipment, etc.)
- 2) Choose grain of the result business process key's individual transaction
- 3) Choose the dimensions
- 4) Choose the measures for each fact table record.

20/12/2023

## Focus of Database Design

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

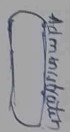
## Three-Tier DW Architecture (8-10 M) \*\*\*



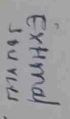
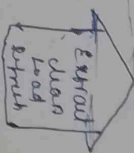
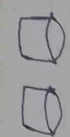
Top-Tier  
front end  
tools



Middle Tier  
OLAP server



Bottom Tier  
Data Mart



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

- four:

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

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

## Enterprise Data 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 slow & low.

### 2) Multidimensional OLAP server

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

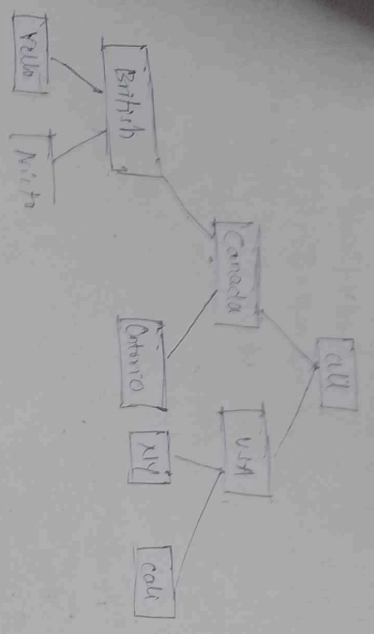
→ Through indexing info. can be retrieved.

### 3) Hybrid OLAP server

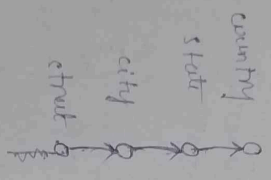
→ It is combination of both relational and multidimensional OLAP servers.



# Concept hierarchy

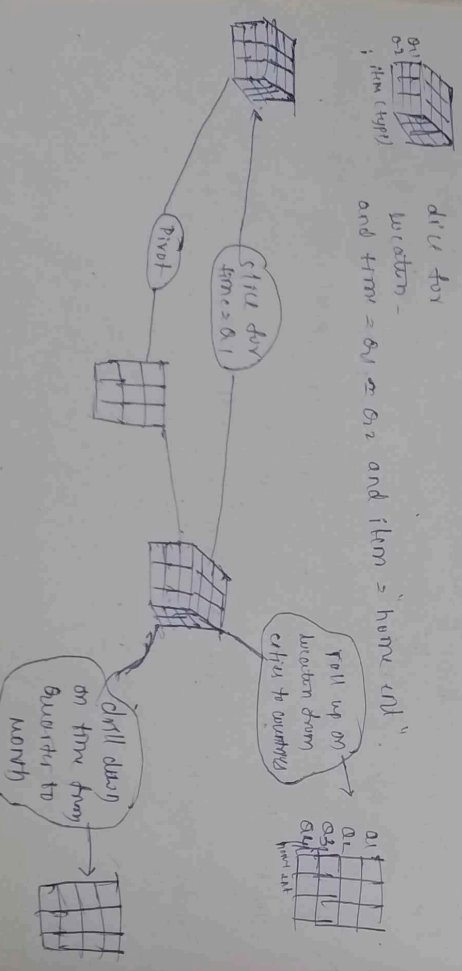


def: Concept hierarchy defines a sequence of mappings from  
 at of low level concepts to higher level. The mappings form  
 a concept hierarchy for the dimension location mapping a  
 set of low level concepts (cities) to higher level concepts  
 (countries)



OLAP operation in multidimensional data model

07/12/2023



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., months, week, day

slice operation: Performs selection on 1D of the given cube.  
 e.g., time = Q1

Pice operation: dividing subquery by performing a selection on two or more dimensions.

Pivot: It provides an alternative presentation of data.  
 e.g., Rows → columns or columns → rows

## Datawarehouse Architecture

→ 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

different views regarding the design of data warehouse

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



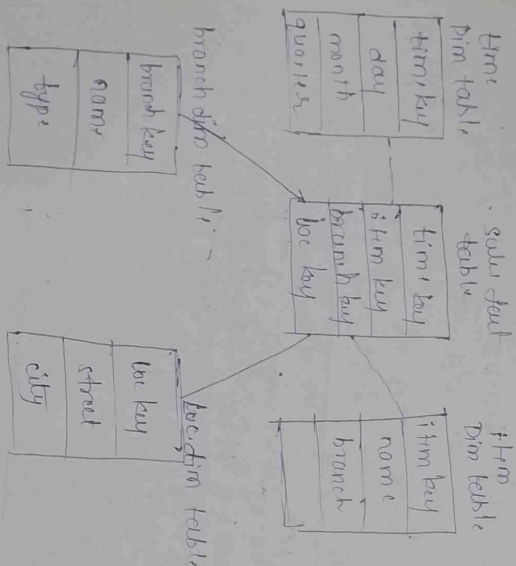
Star, snowflake and fact constellations.

Schema for multi-dimensional DB. (10-15 M\*)

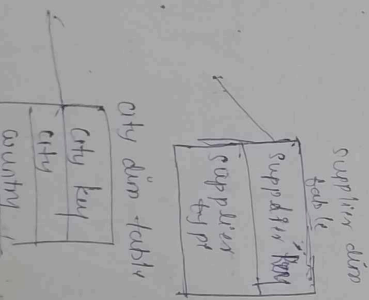
1) star schema of datawarehouse for sales

→ Def: The data warehouse contains a large center table (fact table) containing the bulk of data, with no redundancy.

→ set of smaller dimension tables (Dimension tables) one for each dimension. (Many dim table one fact table)



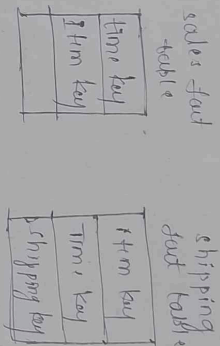
2) Snowflake schema  
It is the variant of the star schema table where some dim tables are normalized then by further splitting the data into additional tables.



3) Fact constellation schema

It requires multiple fact tables to share dimension tables. This kind of schema can be viewed as collection of star schema. It is called as galaxy schema or fact constellation schema.

(Many fact table one dim table)



Examples for defining star, snowflake and fact constellation schema

→ SQL based query language - DML

→ DW defined using 2 languages: primitive

→ one for cube def

→ one for dimension def

→ The fact cube & cube names < dimension list > : (measure list)

→ define dimension < dimension names > as (< attribute list >)

e.g., star schema

define cube sales-stor [time, item, branch, location];

dollars-sold = sum (sales-in-dollars), with-sold = count (\*)

define dimension time as (time-key, day, month, year)

define dimension item as (item-key, name, brand)

define dimension branch as (branch-key, item)

define dimension location as (location-key, street, city)

e.g., snowflake schema

define cube sales-stor [time, item, branch, location, shipping]

define dimension location as (location-key, street, city)

define dimension item as (item-key, name, type, supplier)

define dimension item as (item-key, name, type, supplier)

Star Schema Example:  
Fact Table: sales  
Dimensions: time, item, branch, location  
Query: sum(sales) by time, item, branch, location  
Result: 1000