

Data Mining & Warehousing.

3/2/2023

→ Data, Pattern, Attribute

→ discovery:- database → data mining tools → knowledge

→ Eg of data mining:- FD, M, DW, A, MB.
not structured

→ How data mining is used. (Points)

Point 2:- Excel to digital.

→ Data: Origin of data mining:-

data mining before machine learning / AI algo. * High dimension.

→ Data Mining tasks:- 1) Classification:- Fish goes to water (Eg) (New entity in then define one class)

2) Regression:- Not define to classification.

Eg:- Bird can fly or not (Yes or no only)

3) Clustering:- Eg:- 1000 data then find $K=10$ to every (Many) then find distance from first 10. & lastly small distance we take.

4) Dependence.

5) ...

→ Data Mining Method:- (Not all required (54%))

→ Why data preprocessing?

→ Data Cleaning.

- To classify we have to know the ~~typ~~ category of the items.

Classification techniques or algo:-

Regression, distance, ----

Classification Example:- Grading.

- Partitioning based classification.

Distance based classification.

Classification using Regression → Deviation
→ Prediction

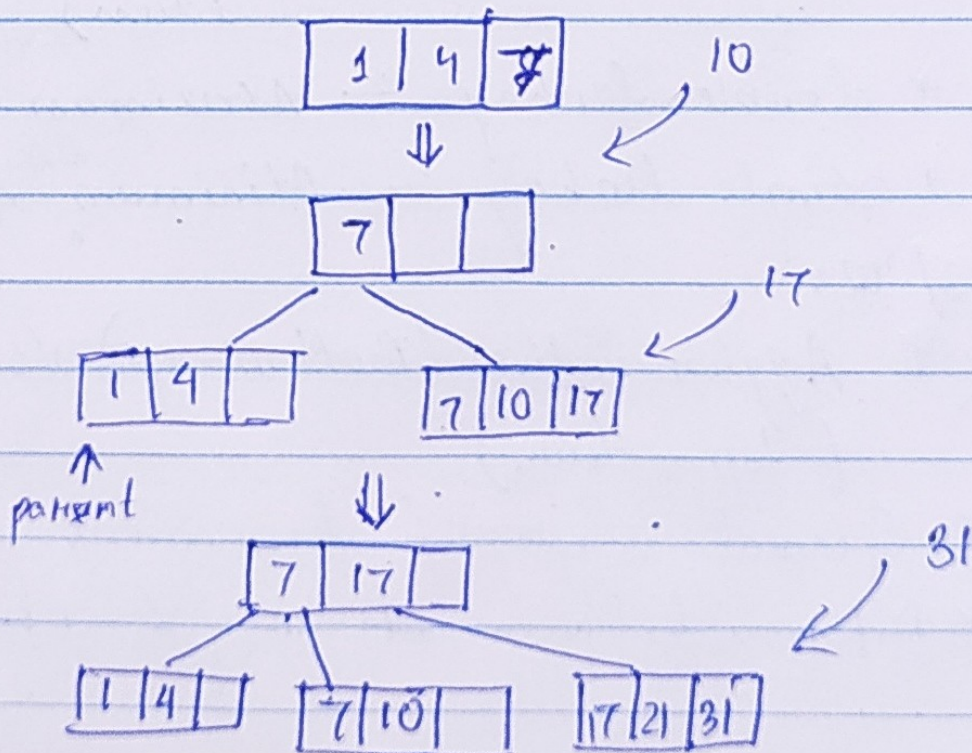
* KNN:- ~~K-Nearest~~ Algorithm

13/3/2023

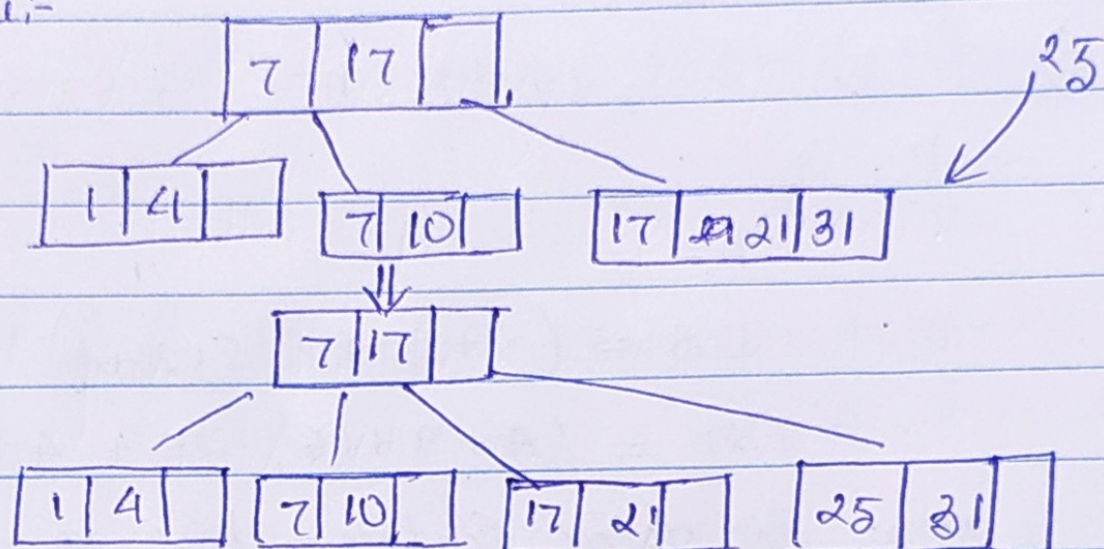
→ BIRCH Algorithm:- (Hierarchical Algorithm)

Eg:- (1, 4, 7, 10, 17, 21, 31, 25, 19, 28, 42)

Soln.



Continue:-



* Phases of BIRCH Algo:- (2-phases)

(a) Phase-1:- (a) Cluster feature tree.

(ii) Scan the database.

~~(b)~~

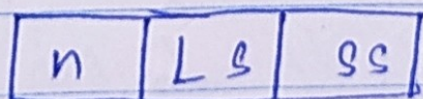
(iii) To build & initialize the cluster.

~~(b)~~

(iv) Leaf node holds many small & tight clusters.

(b) Phase-2:- (i) BIRCH applies a clustering algo to clustered in leaf node of the CF tree.
(cluster feature)

* Every node has three values:-



No. of clusters

↓
Ls
distance

ss Squared distance

$$\Rightarrow (2,5), (3,2), (4,3) \quad n=3, LS=2, SS=?$$

Solⁿ

$$n=3$$

$$LS = (2+3+4, 5+2+3) = (9, 10)$$

$$SS = (4+9+16, 25+4+9) = (29, 38)$$

Now, the cluster Centre = $LS/n = \mu_0$

$$\mu_0 = \frac{(9, 10)}{3} = (3, 3.33)$$

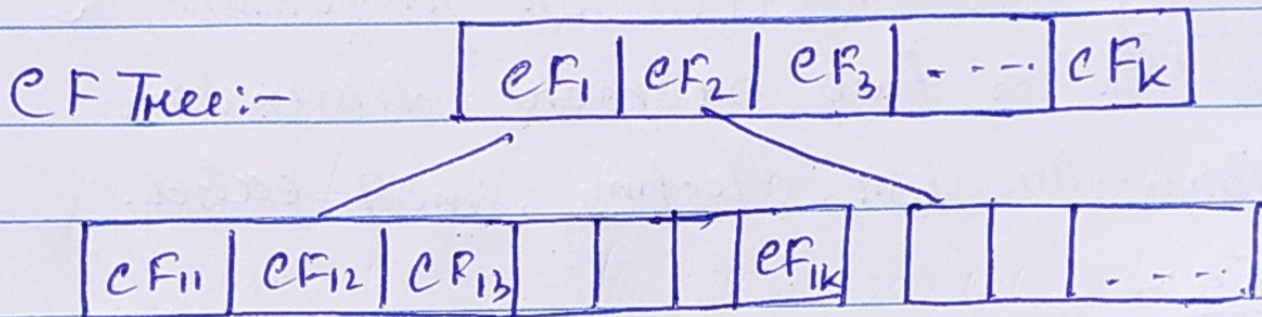
* CF Tree:-

- (a) A height balanced tree, it stores clustering features.
- (b) ~~Long~~ ^{Non} leaf nodes stores sums of the OF's of their children.
- (c) CF tree has two parameters branching factor B and threshold for value T .
where, B = maximum no. of ~~Long~~ ^{non} leaf nodes.
 T = max^m diameter of subclusters at leaf node.

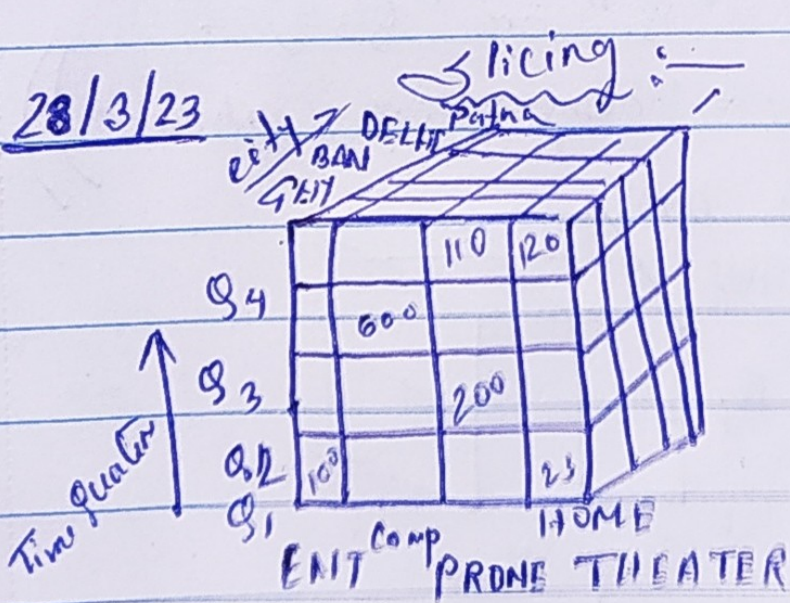
(d) R, D Reflects the typicality of the cluster around the centroid. It is measured as —

$$R = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}} = \sqrt{\frac{nSS - 2LS^2 + nLS}{n^2}}$$

$$D = \sqrt{\frac{\sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2}{n(n-1)}} = \sqrt{\frac{2nSS - 2LS^2}{n(n-1)}}$$



We can find R & D from this CF tree.
 Data structure used in CF tree:- B+ tree.



Drilling (up & down)
 Pivoting.

Sales warehouse → items

Slicing Q3

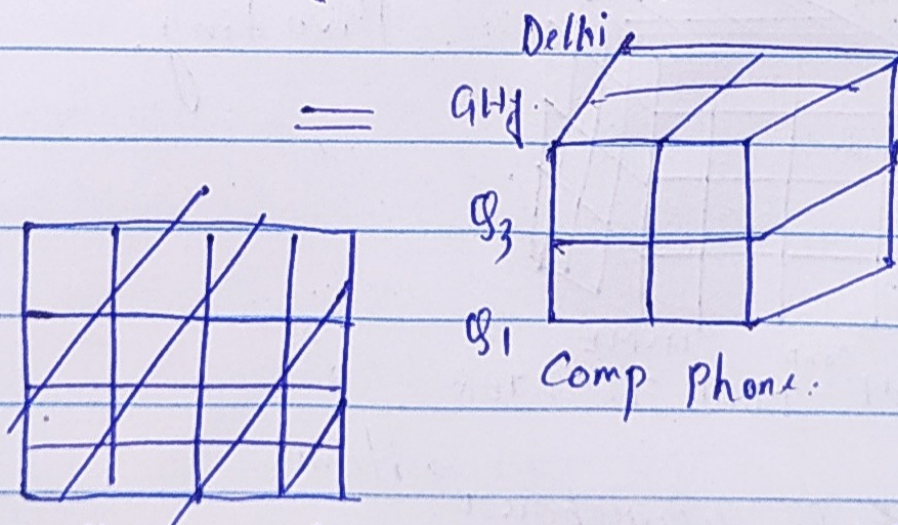
data cube
 Slice
 $C[quality, city, item]$
 time = ' Q_2 '

④ Slicing reduces dimension by one.

→ Dicing:- It is used for performing a selection on two or more dimensions. Basically for selecting small cubes.

Eg:-

DICE FOR ((city = "GHY" OR "DELHI") &
 (Time = " Q_1 " OR " Q_3 ") &
 (ITEM = "comp" OR "phone"))



→ Drilling:- It means moving up & down along classification hierarchical levels.

Drilling up:- It performs a detail to an aggregated level.

Drilling down:- It is reverse of drill up, it navigates from less detail to more detail.

* Drill up/down \Rightarrow Roll up/down.

Eg:- (a) Roll-up location = "Get from Cities to ^{State} Countries"

c [ITEM, QUARTER, CITIES]

##

(b) Roll down

From Quater to month

on
c [ITEM, QUARTER, CITIES]

(c) Roll-up on location (num of phone items sold in '02' in Delhi & GHy)

~~is~~ rotates

(d) Pivot is a rotate is a virtualization operation which rotates the data axes to provide an alternative representation of the same data.

SLICE on Time = "Q1" \in [ITEM, QUARTER, CITIES]

