

ok 2/3  
Monday

4-d-a-1

## Module-1

### Data Mining

Extract of interacting (non-trivial, implicit, permanently unknown and potentially useful) information or patterns from data <sup>in large database.</sup> Processed form of data is known as information.

useful data.



KDD - knowledge discovery in database.

knowledge extraction.  
data/pattern analysis.

data archaeology.

data dredging.

information harvesting.

business intelligence.

No! data mining.

(Deductive) query processing.

KDD is also known as datamining.

### 8 steps in KDD

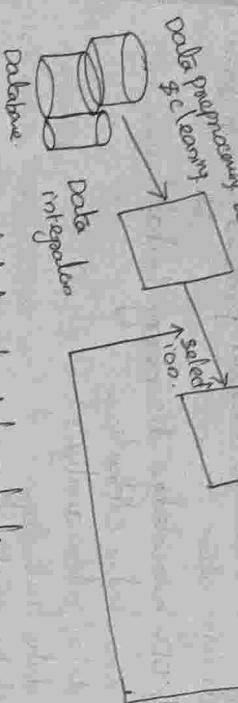
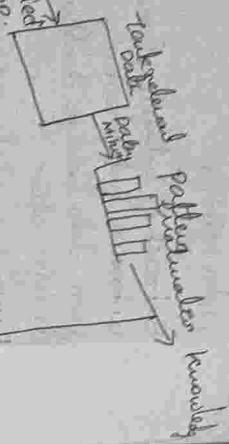
1, Data preprocessing and cleaning.

2, Data selection and Transformation.

3, Data Mining

4, Pattern evaluation.

Data mining: the core of knowledge recovery process.



creating a target data set : data selection.

Data cleaning and preprocessing : removing unwanted data.

2) find useful features/dimensionality

choosing functions of data mining  
summarization, classification, regression, association, clustering

choosing the mining algorithm

Data mining : search for patterns.

Data mining : search for patterns.

Pattern evaluation is knowledge

visualization, transformation, removing redundant patterns.

use of discovered knowledge.

Architecture a Typical Data mining system

6 Major components

1. Database, data warehouse, intranet, or other information repository

2. Database or database known server.

3. Knowledge base.

4. Data mining engine

5. Pattern Evaluation module

6. user interface

1/2/22: Data mining functionalities

Tuesday: Data mining functionalities

• Concept description: characterization and discrimination

• Generalize summarize and constant data:

characteristic eg: dry vs wet regions

Association (cooperation and causality)

Multidimensional vs single dimensional analysis

classification and prediction.

Finding models that distinguish classes or concept for prediction

Prediction

we are grouping into different classes.

Eg: classify cars based on gas mileage.

Prediction: we are predicting new values from previous values.

Regression is also prediction algorithm.

Classification include algorithms, decision trees, classification rule, neural net.

Simple, multi linear, polynomial - analysis

Cluster analysis / clustering

Class label is known: Group data to form a new classes.

Eg: cluster based houses to find delineation pattern.

cluster based on the principle maximize between the intra-

similarity and minimizing the inter-class similarity.

Eg:

1.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

2.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

3.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

4.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

5.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

6.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

7.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

8.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

9.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

10.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

11.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

12.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

13.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

14.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

15.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

16.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

17.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

18.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

19.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

20.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

21.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

22.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

23.  $\text{X}_1, \text{X}_2, \dots, \text{X}_n$  are observations to induction.

## Data warehouse

A data warehouse is subject-oriented, integrated, time-varied, non-volatile collection of data in support of management decision making process.

Decision support database

Data warehousing: The process of conducting and using data warehouse.

Time variant: entirely different from database.

Non-volatile: operational update of data doesn't occur in

the data warehouse environment.

Initial loading of data screen of data

Outlier analysis: A data object that does not comply with general behavior of the data.

Outlier Analysis: It can be considered as noise or exception but it is useful in fraud detection and even analysis

## Trend and evaluation analysis

Trend and deviation: regression analysis similarly based analysis

Segmented path mining, periodicity analysis other pattern: directed on statistical analysis

X. Major issues in data mining

Mining methodology and user interaction Mining different kinds of knowledge in database.

Performance and scalability

Issues relating to the diversity of the data types.

Issues related to application and social impact

Human interaction

Query fitting changing data

Always Data warehousing

Dimensional data

Data cube

Multi-dimensional

Relational

Time

Geographic

Organizational

Product

Customer

Supplier

Market

## Database vs Datawarehouse

Feature	Database	Datawarehouse
Purpose	It is designed to record data.	It is designed to analyze data.
Processing method.	The database uses the online transaction processing (OLTP).	The datawarehouse uses the analytical processing (OLAP).
Usage	Helps to perform random-end operation for your business.	Allows you to analyze data up to date details, historical, summarized flat relational isolated consolidated multidimensional, integrates repetitive read/write/loader/hash lots of scans.
Data	Simple because they are denormalized.	Complex as they are normalized.
Orientation of data	Application-oriented collection A subject-oriented collection of data.	Tablet & join complex as they are normalized.
Designing tools	ER-modelling techniques are used.	Data modelling techniques are used.
Data type	Data stored in the database current and historical data is up to date.	Data stored in data warehouse may not be upto date.
Query type	Simple queries answer complex queries are used.	Complex queries answer simple queries are used.
Data summary	Detailed data stored in the database.	It stores highly summarized data.
OLTP - Online Transaction Processing.		

## OLTP

Users	clerk, IT professionals, knowledge workers.
function	day-to-day operation.
DB design	decision support application oriented.
Access	subject oriented.
Unit of work	short, simple transaction complex query.
# records of access	few millions.
Hours	1000s
DB size	100MB-CB
Index	transaction throughput query, throughput response

## OLAP

13/2/23	it's a multidimensional data model
Monday	A data warehouse is based on a multidimensional which view data in the form of a datacube. Data is represented in the form of data.
item	location="Chicago" loc="Newyork" loc="Toronto" loc="Vancouver"
item	item

Q1  
Q2  
Q3  
Q4

Data model

mekey

data cube representation of data  
Supplier = "Sup1" — ~~Supplier = "Sup2"~~

**13223** A multilane bridge is based on a multi-lane bridge. The term often used

A data warning is given in the form of a table which views data in the form of data and more.

17. *Lab* is rep. at *Montreal* loc = "Montreal" loc = "New York" loc = "Toronto"

Supplier = "SUP3'

home home end. comp. phones) end. comp.

time and corp. phone. ~~Heik.~~ 201.

$Q_4$   $Q_3$   $Q_2$   $Q_1$

### Data cube

Data cube allows data to be modelled and viewing multiple dimensions. It can be defined by dimension.

and Gets  
Dimensions

Dimensions are the perspectives or entities with which a table, store exists.

A lamp of Semur  
home phone


## Example of dimension table.

time
time-key
day
day-of-the-week
month
quarter
year

branch
branch-key
branch-name
branch-type

eg of fact table

time-key.
item-key
branch-key.
location-key.
units-sold.
dollars-sold
avg-sales.

Problem

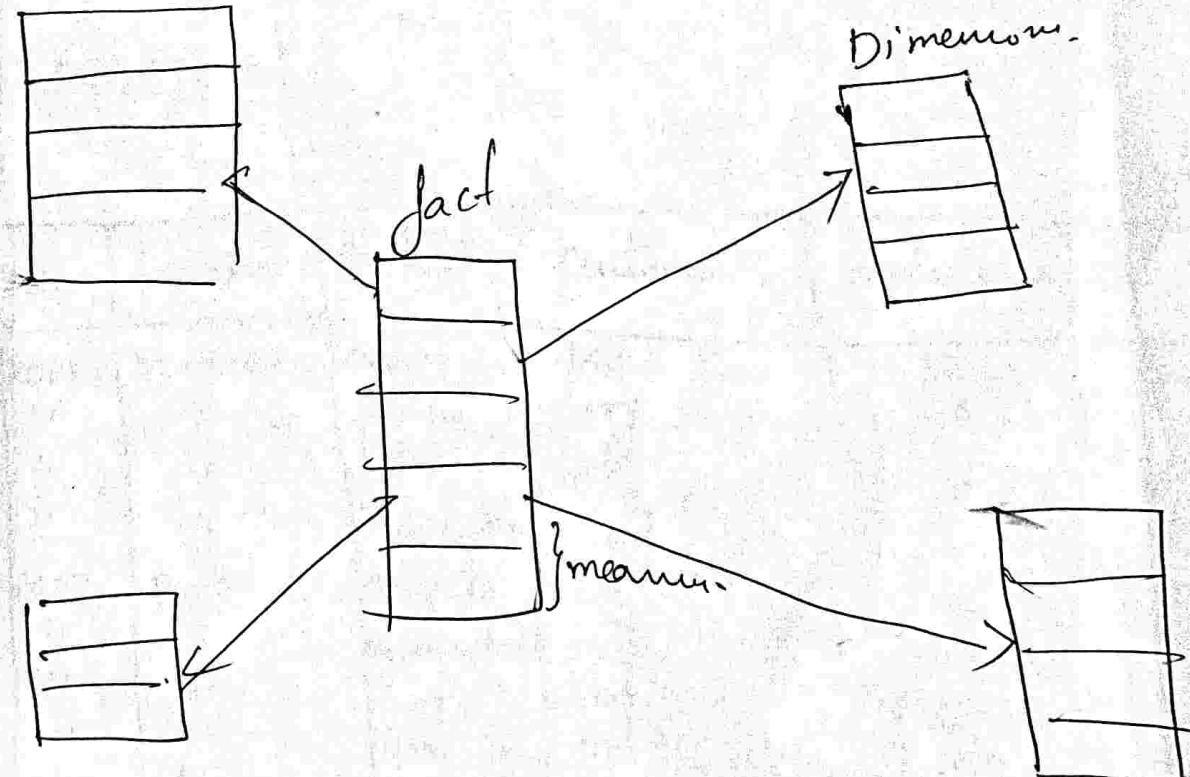
measures

## X. Conceptual Modelling of Data warehouse / Data warehouse schema

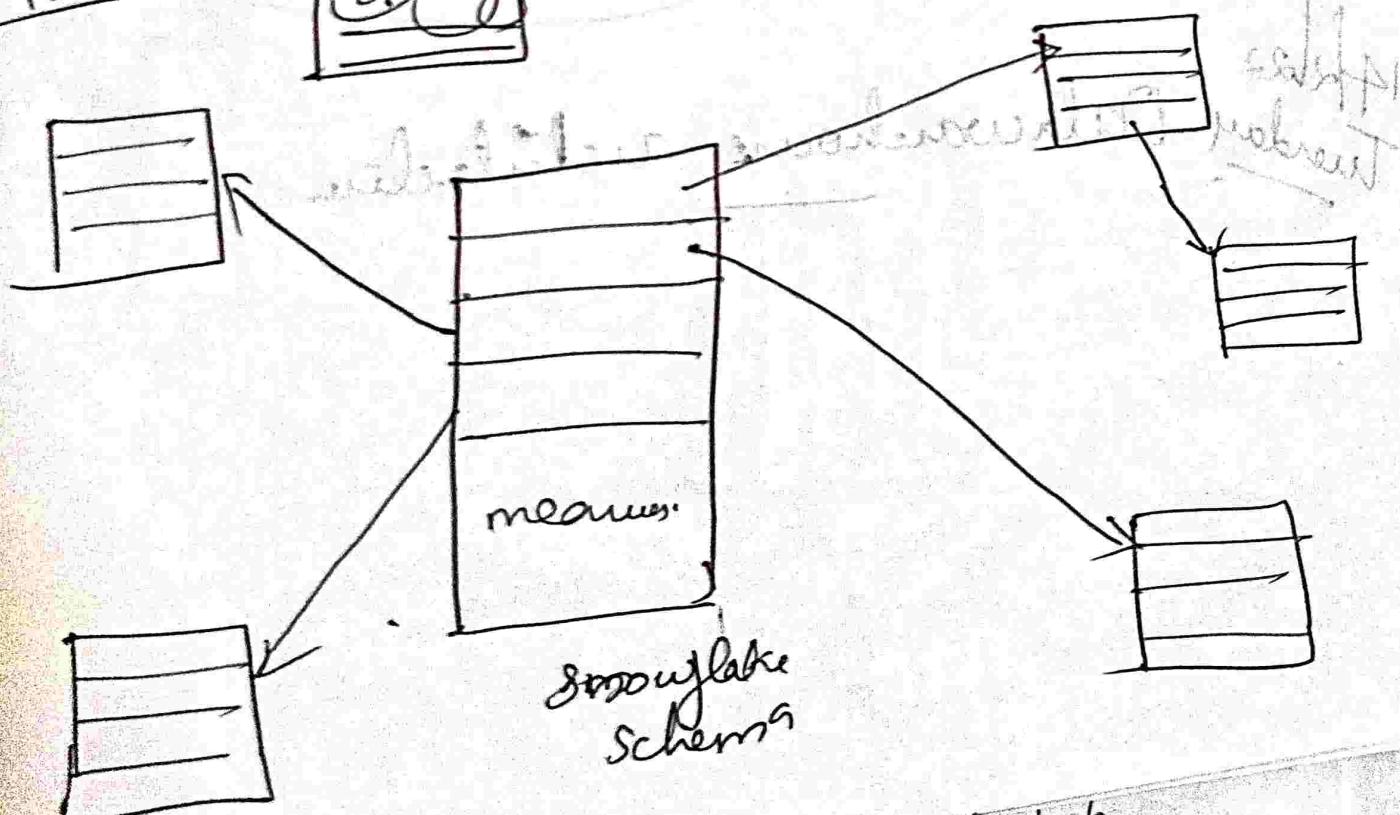
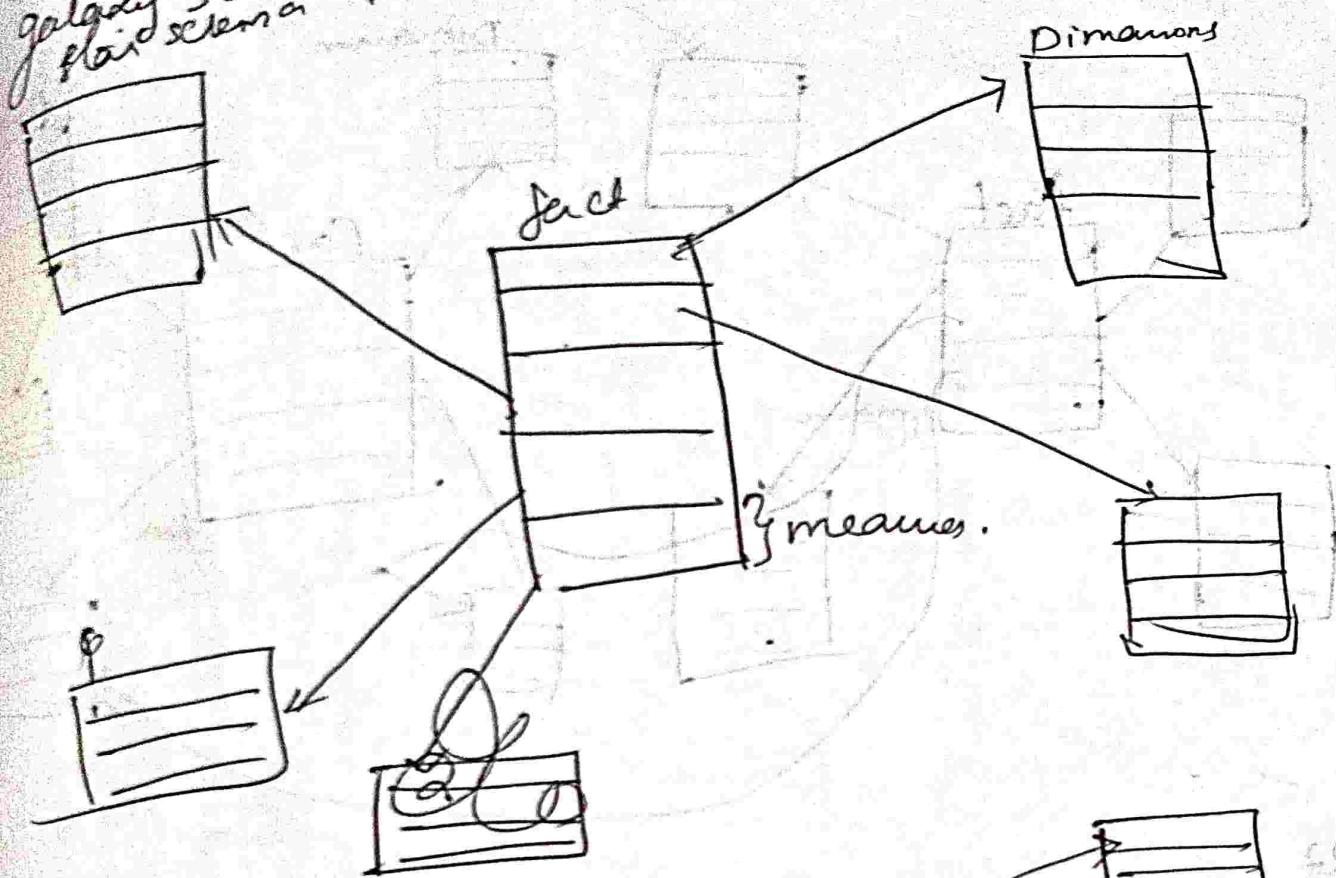
star schema - fact table <sup>connected</sup> is the middle to a set of dimension table.

Snowflake schema : A refinement of star schema to normalize into a set of smaller dimension tables, forming a shape similar to.

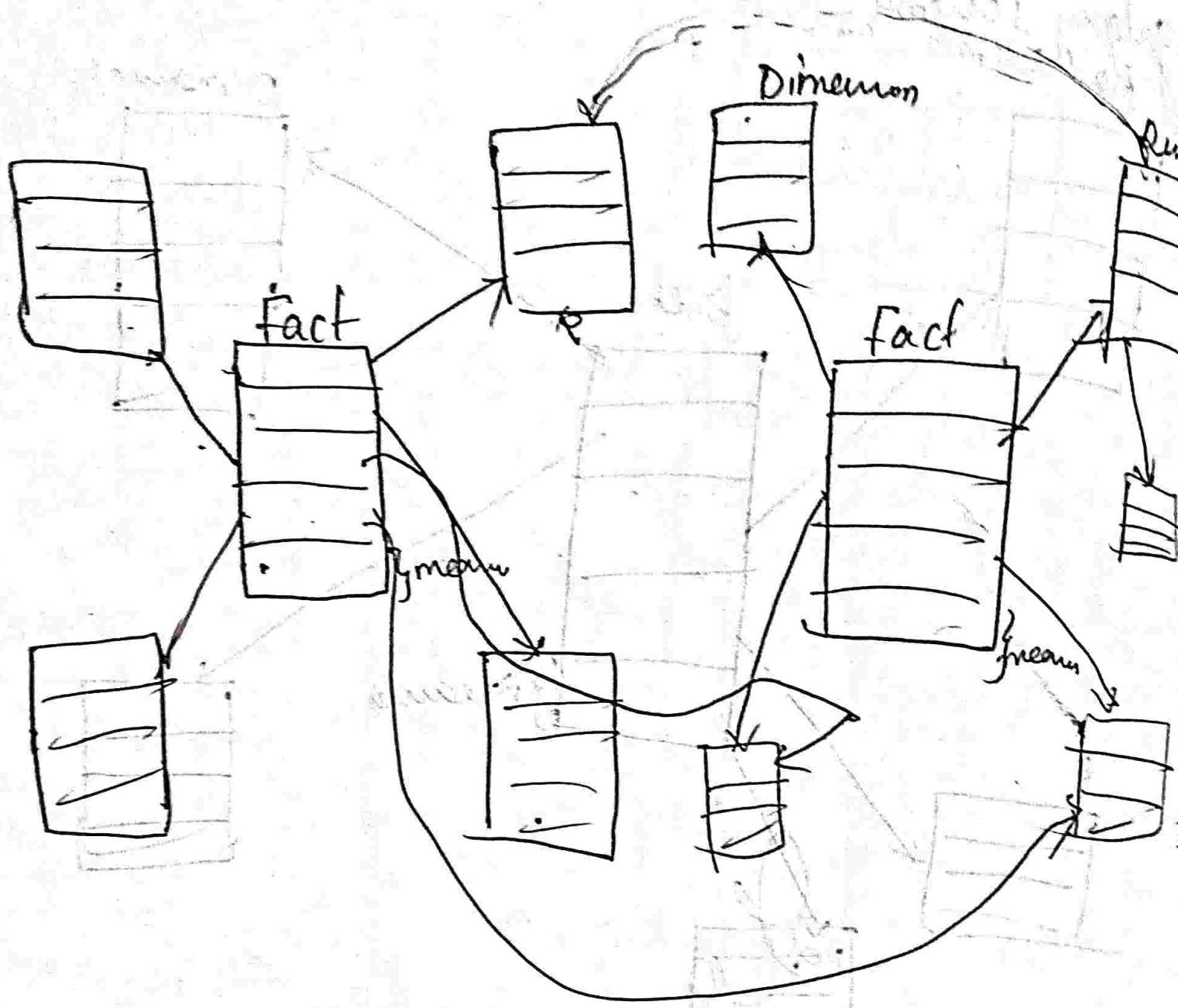
Galaxy schema : Multiple fact tables share dimension tables, viewed as a collection of states through states therefore called galaxy schema .



multiple fact tables share dimension tables  
viewed as a collection of stars therefore called  
galaxy schema.  
flat schema.



# Galaxy schema



14/2/22  
day Datawarehouse architecture

14/12/23  
Tuesday

## Data ware house architecture

Material  
architecture

Data mining  
application

Operational  
Data



Raw



Cleaning,  
integrating  
process



Background process

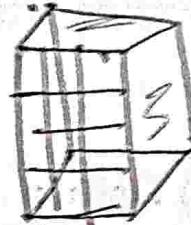
METADATA

Reconciled data.



ODS.

Derived  
data



DW  
service

Analytic



Time

Data  
marts



DW

marts

DW

used to document the dimension hierarchy by one down.

2011	Hydabad	Music-system	Ind	22
"	"	compute	Org	15
"	"	Ent	Ind	3

Ind - Individual  
Org - Organization

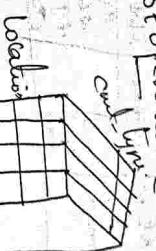
### 1) Roll-up operation

Increase the hierarchy by one dimension up. i.e.,  
used for summarization of data dimension.  
e.g.: product (Home, entertainment, computer system)  
the we are categorizing product attribute into 3 categories  
Home appliances, entertainment, computer systems after rolling up  
operation like table will be.

Table 2

Year	Location	Product-type	Customer	Total sales
2008	Mumbai	Home-appliances	Ind	32
"		"	Org	34
"		entertainment	Org	12
"	Pune	Music-system	Ind	10
2009		Computer-system	Ind	43
		entertainment	Ind	45
2010	Chennai	Home-appliances	Org	09
"		entertainment	Ind	26
"		Computer-system	Org	13
2011	Hyderabad	entertainment	Ind	29
"		Computer-system	Org	65
"	H A		Org	63

### 3) Pivot operation (rotate)



This is used to view data cubes in different form.

view 3D



2D planes

4) slice : It is used to extract a slice of the original cube corresponding to a single value of a given dimension. Similar to select operation in RDBMS.

for e.g. Slice Table 2 to on location = "Hyderabad".

Year	Location	Product-type	Customer	Total sales
2011	Hyderabad	entertainment	Ind	29
"	"	Computer-system	Org	63

### 2) Roll-down / Drill down

Reverse operation of roll-up.

### 5) Dice operation

It is used for selecting two or more dimensions and one more value.

e.g. "mumbai" ("hyderabad" or "chennai" and product type)

(entertainment or H.A.)

After applying dice operation from Table 2.

year	location	product	customer	Total sales
2010	Chennai	H.A	Org	89
"	Hyderabad	Entertainment	Ind.	26.
2011	"	"	Ind.	29
"	"	H.A	Org	03.

### 6) Drill Aways

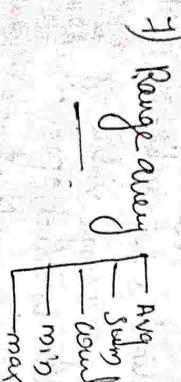
Analyze the cells of a datacube, using same set of dimensions.

Used to move data from a datacube to another.

Applied only when data cubes are having same/parent dimension.

e.g. sales

Prod	
	Prod..



It applies a group aggregation operation over a set of selected cells.

For eg. find the total no. of sales for a period 2008-2010. Thus the range query will return the value by aggregating the values, total sales from 2008 - 2010. It can be applied on dimensions with

numerical values.  
Range sum query (sum).

Year	Location	Product	Customer	Avg sales per month
2008	Mumbai	Orange juice	Ind	10
"	"	"	Org	03
"	"	Tv	Ind	07
"	"	"	Org	01

## Module - 2

### Data Preprocessing

3 types of data

Incomplete

Noisy

Inconsistent

why data preprocessing?

real world data

incomplete - redundancy of data.

incomplete - redundant datarow.

Noisy - unwanted data row.

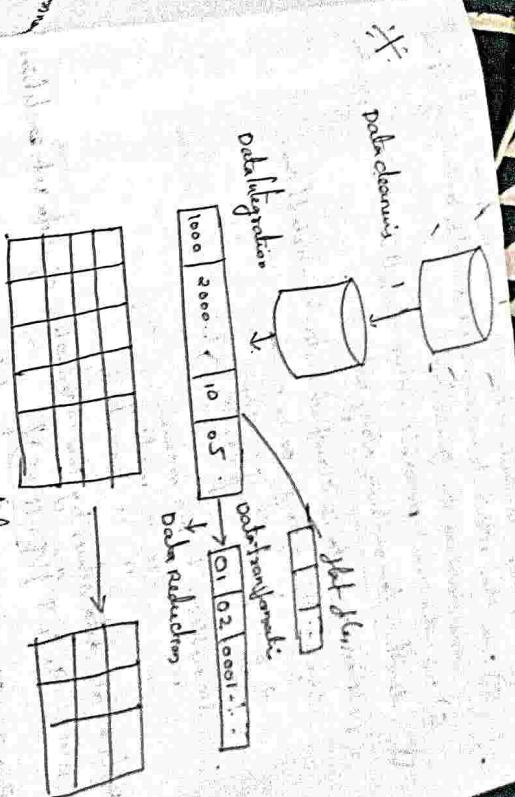
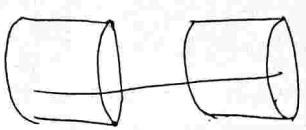
Inconsistent

Naming confusion : same data and different column

Normalization

+ Major steps in data preprocessing

- (i) Data cleaning
- (ii) Data integration
- (iii) Data transformation
- (iv) Data reduction



Data cleaning : remove noisy data  
filling missing

Data transformation : To convert into common format  
(normalization)

Data reduction : To reduce its size after data transformation.

23/2/23 - Wednesday: Data cleaning : Filling the missing value, smoothing noisy data  
identifying or removing outliers, dealing inconsistency

I Handling missing values

1) Ignore the tuple.

2) Method is not very effective

2) Difficult in the missing values manually. May not be possible.

This approach is time consuming. May not be feasible.  
3) use a global constant to fill in the miss value.

- 4) Use the attribute mean to fill in the missing values  
 5) are the attribute mean for all samples belong to it  
 1 same class as the given tuple.

- c) use the most probable value to fill in the missing values

## II Handling/Smoothing Noisy Data

(i) Binning: Data is classified into buckets/bins.

Consider the data for price

$$\text{Bin 1: } \begin{array}{|c|c|} \hline 8, 15 & 21, 21, 24 \\ \hline \end{array} \quad \text{Bin 2: } \begin{array}{|c|c|} \hline 25, 28, 24 & 25, 28, 24 \\ \hline \end{array}$$

partitions

smoothing by bin mean.

$$\begin{array}{l} " \\ " \\ " \\ " \end{array} \quad \begin{array}{l} \text{median} \\ \text{boundaries} \end{array}$$

Replace all data values by the mean value of data set each.

$$\text{Bin 1: } 9, 9, 9 \left\{ \frac{(4+15+8)}{3} \right\} = 9.$$

$$\text{Bin 2: } 21, 21, 24 = 22.$$

$$25, 28, 24$$

$$\text{Bin 3: } 25, 28, 34.$$

$$29, 29, 29.$$

smoothing by median.

$$\text{Bin 1: } 8$$

$$\text{Bin 2: } 21$$

$$\text{Bin 3: } 28.$$

$$\begin{array}{|c|c|} \hline 1 : 8, 8, 8 & 2 : 21, 21, 21 \\ \hline 3 : 28, 28, 28 & \end{array}$$

smoothing by bin boundaries.

$$9, 4, 15$$

$$21, 21, 29$$

$$25, 28, 34$$

Apply smoothing technique by min, mean, median, boundaries or  
 the following data

$$13, 4, 12, 4, 15, 13, 9, 11, 12, 3, 4, 12, 14, 23, 1, 8, 9, 0$$

binsize: 5.

sort

$$\begin{array}{|c|c|} \hline 12, 3, 7, 8 & 9, 11, 13, 14, 23 \\ \hline \end{array} \quad \begin{array}{|c|c|} \hline 0, 4, 12, 6, 30, 34, 50 \\ \hline \end{array}$$

$$\begin{array}{l} \text{Bin 1: } 1, 2, 3, 4, 8 \\ \text{Bin 2: } 9, 11, 13, 14, 23 \\ \text{Bin 3: } 24, 26, 30, 34, 50. \end{array}$$

smoothing by bin median.

$$\begin{array}{l} \text{Bin 1: } 4, 2, 4, 2, 4, 2 \\ \text{Bin 2: } 14, 6, 14, 6, 14, 6 \\ \text{Bin 3: } 32, 8, 32, 8, 32, 8 \end{array}$$

smoothing by bin boundaries.

$$\begin{array}{l} \text{Bin 1: } 3, 3, 3, 3, 3 \\ \text{Bin 2: } 13, 13, 13, 13, 13 \\ \text{Bin 3: } 30, 30, 30, 30, 30. \end{array}$$

Regression (prediction)

If it is also used for noise recovery.

Type of regression:  
 Simple linear.

Multiple  
 Polynomial  
 Radial basis function

$$\begin{array}{|c|c|} \hline 1 : 8, 8, 8 & 2 : 21, 21, 21 \\ \hline 3 : 28, 28, 28 & \end{array}$$

Data cleaning as a process.

Saturday step 1: Discrepancy detection.

2: Data transformation.

Discrepancy

Human error in data entry

## Avoid discrepancy

During migration

Commercial tools.

Data scrubbing tool.

Data auditing tool.

Data transformation

Data Migration tool

ETL (Extraction/Transformation/Loading).

clustering

help to avoid errors in schema entity relationship problem.  
Data redundancy (repetitions of data)  
using normalization to avoid redundancy. Derived attribute  
are also examples of redundancies.  
solution to data redundancy

Correlation Analysis

For numerical attributes we can evaluate the correlation b/w  
attributes A & B known as pearson's product moment coefficient  
(named after its inventer karl pearson)

$$\gamma_{A,B} = \frac{1}{N} \sum_{i=1}^N (a_i - \bar{A})(b_i - \bar{B})$$

$N \sigma_A \sigma_B$

$$\text{Assume} \quad (\sum_{i=1}^N (a_i b_i)) = N \bar{A} \bar{B}$$

No. of tuples

$\bar{A}$  - No. of tuples.

$\bar{B}$  - respective mean values of A & B in tuples.

$\sigma_A$  &  $\sigma_B$  - respective standard deviation.

A, B - attributes

$\gamma_{A,B}$  - correlation coefficient.

$-1 \leq \gamma_{A,B} \leq 1$ .

If  $\gamma_{A,B} > 0$ , A & B are very correlated. which mean the value

of A increase as the values of B increase.

Highly correlated if  $\gamma_{A,B} = 0$ , A & B are independent there is no correlation b/w them.

## II

### Data Integration

The data integration is the process of combines data from multiple sources into a coherent data store.

Issues in data integration

- Entity identification problem
- Data Redundancy problems.

- Resolution of data value conflicts
- Defects and resolution of data value conflicts

- Entity identification problem

Eg: Data analyst or compute cannot ensure that customer id in one database and cust-number in another refer to the same attribute, even though

if  $r_{A,B} < 0$

A & B are negatively correlated, where the values of one increase, the values of the other attribute decrease.

If  $r_{A,B} \neq 0$  A & B dependent to each other.

### Correlation analysis for categorical attribute

A correlation relationship b/w 2 attributes A, B can be discovered by a  $\chi^2$  (chi-square) test. Suppose A has c distinct values, namely  $a_1, a_2, \dots, a_c$ . B has r distinct values, namely  $b_1, b_2, \dots, b_r$ .

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

The data tuples described by A & B can be shown as a contingency table, with the C value of A

$O_{ij}$  - observed frequency (actual count) of the joint event  $(A_i, B_j)$

$E_{ij}$  - expected frequency of  $(A_i, B_j)$

The sum is computed over all of the

$$E_{ij} = \text{count}(A=a_i) \times \text{count}(B=b_j)$$

### Problem of categorical analysis

- Suppose that a group of 1,500 people was surveyed. The gender of each person was noted. Each person was polled as to whether their preferred type of reading material was fiction or non-fiction. Thus, we have 2 attributes, gender and preferred reading. The observed frequency (or count) of each possible joint events is summarized in the Contingency, where the numbers in parentheses are the expected frequencies.

	Male		Female		
$b_1$ Fiction	250 (1,1)		200 (2,1)		450
$b_2$ Nonfiction	50 (1,2)		1000 (2,2)		1050
	300		1200		1500

Hypothesis

Gender and preferred reading are independent each other.

$B_{ij}$  = Total cost

$$e_{ij} = \frac{(a_{ij} \times b_{ij})}{N}$$

$$e_{11} = \frac{a_{11} \times b_{11}}{N}$$

$$= \frac{300 \times 450}{1500} = \underline{\underline{90}}$$

$$e_{21} = \frac{120 \times 300}{1500} = \underline{\underline{360}}$$

$$e_{12} = \frac{a_{12} \times b_{12}}{N}$$

$$= \frac{300 \times 1050}{1500} = \underline{\underline{210}}$$

$$e_{22} = \frac{a_{22} \times b_{22}}{N} = \frac{120 \times 1050}{1500} = \underline{\underline{840}}$$

	$a_1$	$a_2$	Total
	Male	Female	Count
$b_1$ Fiction	250 (1,1) (90)	200 (2,1) (360)	450
$b_2$ Nonfiction	50 (1,2) (210)	1000 (2,2) (840)	1050

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$= \frac{(250-90)^2}{90} + \frac{(200-360)^2}{360} + \frac{(50-210)^2}{210} +$$

$$\frac{(1000-840)^2}{840}$$

$$= 284.44 + 71.11 + [21.90 + 30.48]$$

$$= \underline{507.93}$$

$$\text{degrees of freedom} = (r-1)(c-1)$$

$$= (2-1)(2-1) = 1.$$

If the critical value from the chi-square distribution table is checked against the calculated chi-square value, if calculated value, hypothesis can be accepted. Otherwise hypothesis rejected.

Critical value corresponding to  $1d((r-1)(c-1))$  with maximum significant levels of 0.001 is 10.828.

$507.93 > 10.828$  hypothesis can be rejected. i.e., A & B are strongly correlated, ie, preferred reading and gender are correlated.

Detection and resolution of data value conflicts.

To avoid this attribute functional dependencies and referential constraints in the source system match those in the target system.

28/01/23  
Tuesday

Data transformation

The data are transformed into forms appropriate for mining.

in. i.e., includes binning, minmax, clustering.

28/2/25 Data transformation -  
The data are transformed into form appropriate for mining.

Smoothing is used to remove noise in the data. Such techniques include binning, regression, clustering.

Aggregation: where summary or aggregate are applied to the data.

Eg: the daily sales data may be aggregated so as to compute monthly and annual totals amount. This step is typically used in concluding a data cube for analysis of the data @ multiple granularities.

Generalization of the data, where low-level data are replaced by higher level concepts through the use of concept hierarchies.

Eg: Categorical attributes, like state, can be generalized to higher level concepts like City or Country.

Min-max normalization:  
where the attribute data are scaled so as to fall within a small specified range.

Methods for data normalization:

1. Min max normalization  
2. Z-score

Min-max normalization:  
perform a linear transformation on the original data.

Let  $\min_A$  and  $\max_A$  are the minimum and maximum values of an attribute A.

Maps a value,  $v$ , of A to  $v'$  in the range  $[\text{new\_min}_A, \text{new\_max}_A]$ .

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new\_max}_A - \text{new\_min}_A) + \text{new\_min}_A$$

Eg: we may be wish to add the attribute area based on the attribute height & width. This gives a set of attributes to help the mining process.

Techniques include monitoring, management and training.

<sup>out</sup>  
R<sub>score</sub> or Z-score normalization  
100 : 100 data 99200.

It will encounter the task of reorganization full ahead if we care for it - now or at

Overhead for eq: cancer  
 $\sqrt{99200 - 12000}$

Q: Suppose that the minima and max values for income are \$12,000 & \$28,000.

Medicaid is a program designed like the many "income plus"

respectively. The range, i.e., by min-max normalization, the range [0.0, 1.0], of income is transformed to a value of \$ 73,600 for income, or <sup>estimated</sup> estimated.

$$\sqrt{V} = \sqrt{-m\dot{\phi}^2} \quad (\text{new mass} - \text{new rotation})$$

WILSON - 1940

dīśodātālo

$$V = V - \bar{A}$$

$$\begin{aligned} \min_n &= 12000 \\ \max_n &= 28000 \\ \text{newmin}_n &= 0.0 \\ \text{newmax}_n &= 1.0. \end{aligned}$$

Suppose that the mean  $\bar{x}$  s.d. of the values for the

$$\sqrt{13600 - 12000} \left(1.0 - 0.0\right) + 0.0$$

9.8000 - 12000

$$\cancel{0.160} \cancel{0.01} x + 0.1 = 0.2416$$

11  
11  
60  
88

$$\sigma_A = 16,000$$

Suppose that the mean & s.d. of the values for the attribute income are \$154,000 & \$165,060 respectively. With Z-score normalization a value of -13.600 for income is transformed.

Suppose that the mean & s.d. of the values for the attribute income are \$154,000 & \$165,000 respectively. With z-score normalization, a value of -13.600 for

$$\sqrt{V} = \sqrt{V - A}$$

$$= \frac{73600 - 54000}{16000}$$

$$= \frac{19600}{16000} = 1.225$$

X Normalization by decimal scaling

Normalized by moving the decimal point of values of attribute A.

The no. of decimal points moved depends on max absolute value of A.

A value of  $V_i$  of A is normalized to  $V'_i$  by computing:

$$V'_i = V_i$$

$$= \frac{V_i}{10^j},$$

where  $j$  is the smallest integer such that

$\text{Max}(V_i) < 1.$

Q. Suppose that the recorded values of A are from -986 to 917. The max absolute value of A is 986. To normalize by decimal scaling, we therefore

propriet form.

1. divide each value by  $1,000$  ( $i, j = 3$ ) so that -986 normalizes to -0.986 and 917 normalizes to 0.917. + the following data

$$V' =$$

13/23 use the a method to normalize data

200, 300, 400, 600, 1000  
normalize by setting min=0

(A) Min-max

$$\text{max} = 1,$$

Z-score normalization.

$$(B) V' = \frac{V - \text{min}_A}{\text{max}_A - \text{min}_A} [ \text{new}_{\text{max}} - \text{new}_{\text{min}} ] +$$

$$\text{new}_{\text{max}} = 1.225, 0 = \text{new}_{\text{min}}$$

$$\text{new}_{\text{min}} = 0.$$

$$V = 200$$

$$\text{min}_A = 200, \text{max}_A = 1000.$$

$$V' = \frac{200 - 200}{1000 - 200} [1 - 0] + 0 = 0$$

$$V = 300$$

$$V' = \frac{300 - 200}{1000 - 200} [1 - 0] + 0 = \frac{100}{800} = 10.125$$

$$V = \frac{1}{1000-200} [1-0] + 0$$

$$= \frac{400-200}{800} = -\frac{200}{800} = +0.25$$

$$V = \frac{V-200}{1000-200} [1-0] + 0$$

$$= \frac{100-200}{1000-200}$$

$$V = \frac{V-200}{1000-200} [1-0] + 0$$

$$\sigma^2_{\text{variance}} = \frac{\sum (x - \bar{x})^2}{N}$$

$$SD = \sqrt{\frac{\sum (x - \bar{x})^2}{N}}$$

$$\sigma^2_{\text{variance}} = \frac{1}{3} [(200-300)^2 + (300-500)^2 + (400-500)^2 + (500-500)^2]$$

$$= \frac{800-200}{800} = -\frac{600}{800} = -0.75$$

$$V = \frac{1000-200}{1000-200} = 200-1000 = -\frac{800}{800} = -1$$

$$\sigma^2_{\text{variance}} = \frac{1}{3} [(200-500)^2 + (300-500)^2 + (400-500)^2 + (500-500)^2]$$

$$= \frac{1}{3} [90000 + 40000 + 10000 + 10000] = 80000$$

$$\text{Normalized data} = 0, 0.125, 0.25, 0.35, 1$$

b) Z-score normalization.

$$0.05 = \frac{V}{\sigma_n}$$

$$SD = \sqrt{80000} = 282.84$$

$$V = \frac{200-500}{282.84} = -1.66$$

$$\bar{A} = \frac{200+300+400+600+1000}{5}$$

$$= \frac{2500}{5} = 500$$

$$V = \frac{400-500}{282.84} = -0.35$$

$$V' = \frac{600 - 500}{282.84} = \underline{\underline{0.35}}$$

$$V' = \frac{1000 - 500}{282.84} = \underline{\underline{1.76}}$$

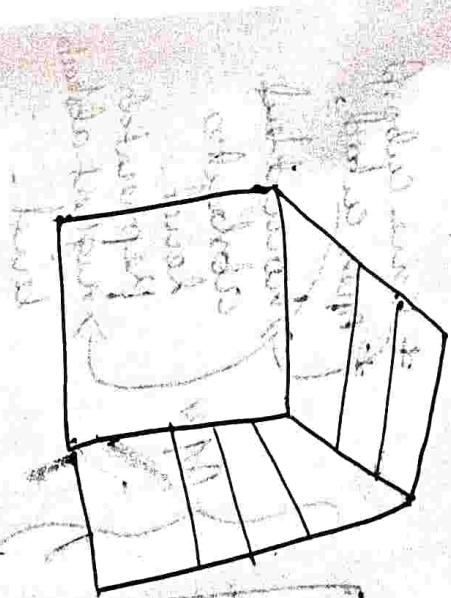
Data Reduction: The process of retaining relevant data by summarizing and eliminating irrelevant attributes from the data set.

Mechanisms used for data reduction:

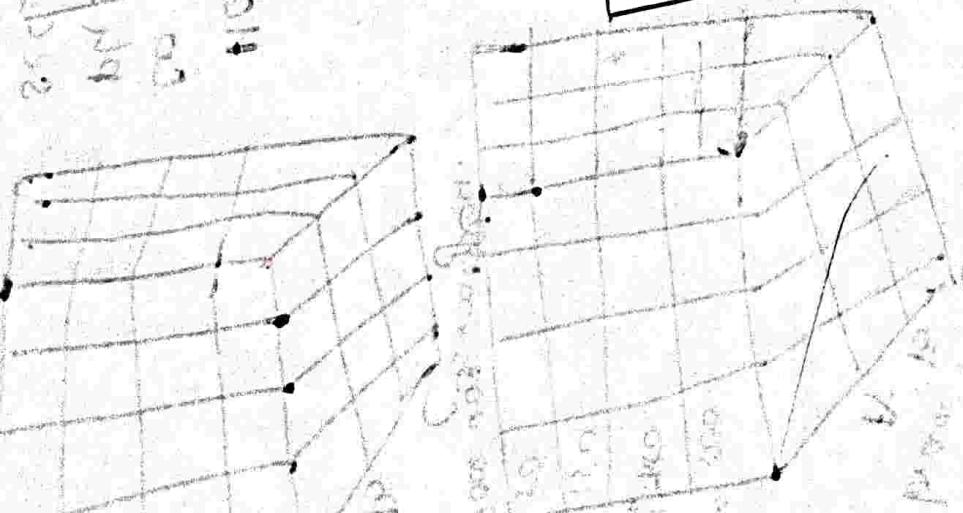
- + Data cube Aggregation
- + Attribute subset selection
- + Dimensionality reduction
- + Numerosity reduction
- + Discretization & Concept hierarchy generation

### (1) Data cube Aggregation

Year	Sales
2002	500



Year 2002	
Quarter	Sales
Q1	500
Q2	500
Q3	500
Q4	500



techniques:

aggregation  
clustering

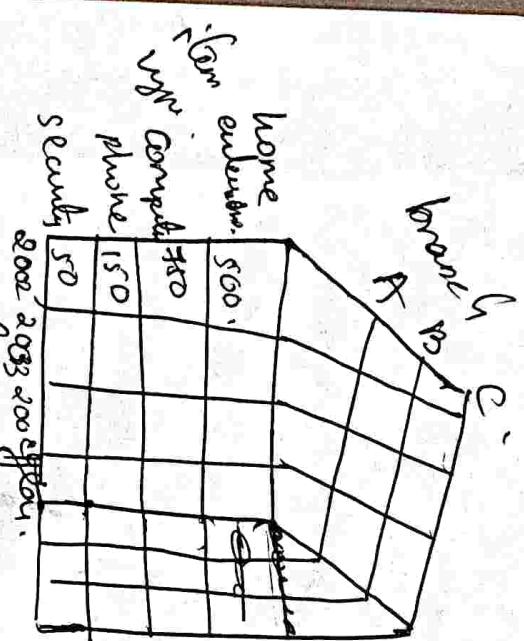
smooth to the  
data such

Aggregation  
data  
eg: the data  
month  
continues

Generates

Year	Sales
2002	100
2003	100
2004	100
2005	0

Year	Sales
2002	100
2003	100
2004	100

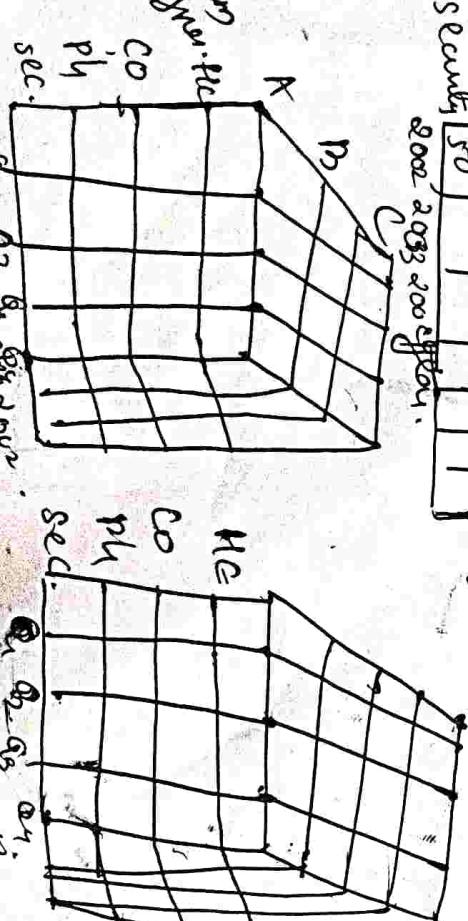
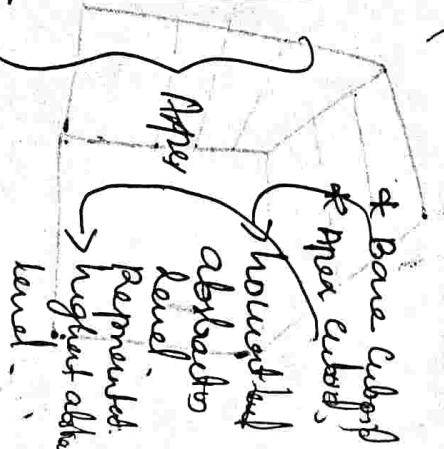


Home extension  
item  
100  
100  
100  
100

Computer  
phone  
150  
150

Security  
50  
50

2002 2003 2004



sec.

CO

phy

sec.

CO

phy

sec.

CO

phy

sec.

CO

phy

sec.

Techniques include binning, maximum clustering, and tree clustering.

Aggregation  
Dimensionality Reduction

6/3/23: If Dimensionality Reduction  
Data: Monday  
eg: the amount  
count  
DWT (Discrete wavelet transforms)  
Original data is signal processing  
Replaced by discrete wavelet coefficients.  
Histogram: If the original data can be  
represented in a more reduced form.

Q. 8 - Dimensionality Reduction are applied in Aggregation & Data Mining.

J. f. Dimensionality Reduction are applied to the DUT (Discrete wavelet transform) signals are transformed into 8 signal frequency components. Replaced by higher quality efficient. However: If the original data can be reduced. From the selected DIP - Then it is known as local encoding.

Numerosity Redn

- to reduce the no. of cows.
- (i) parametric methods
  - Regression.
  - log linear models  $\rightarrow$  probability concept based on parametric.
  - strictly non-parametric models.

- Regression models
- Non linear models
- Non parametric methods
- Histograms

• Club  
Sampling

Histogram: Convert the data to a histogram.

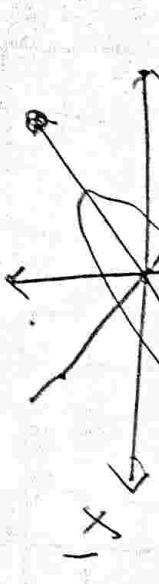
In the case of dimensional data set we have to calculate  $k$  (minimum) principal components where  $k < n$ .  
 They are orthogonal vector  $\perp$  to original data.  
 Principal components are unit vectors calculated from original data  $\therefore p$ .

They are orthogonal vector  $\perp$  to original data.  
Principal components are unit vector calculate from  
original data  $i/p$ .

27

It can be diagrammatically represented.

Y<sub>2</sub> =  $\frac{1}{2}$   
Y<sub>1</sub> = First principle up to Company  
Beechwood  
present.





### Process of partitioning histogram

(1) Equal widths.

→ width of each bucket is uniform.

(2) Equal frequency.

frequency of each bucket is a constant  
each bucket contains same no. of elements

3) V-optimal → Bucket with least variance.

4) Max Diff - Diff b/w each pair of adj values.

7/3/23 Sampling

Tuesday

1) simple random sampling without replacement (SRSWR)

Simple Random Sampling with Replacing (SRSWR)

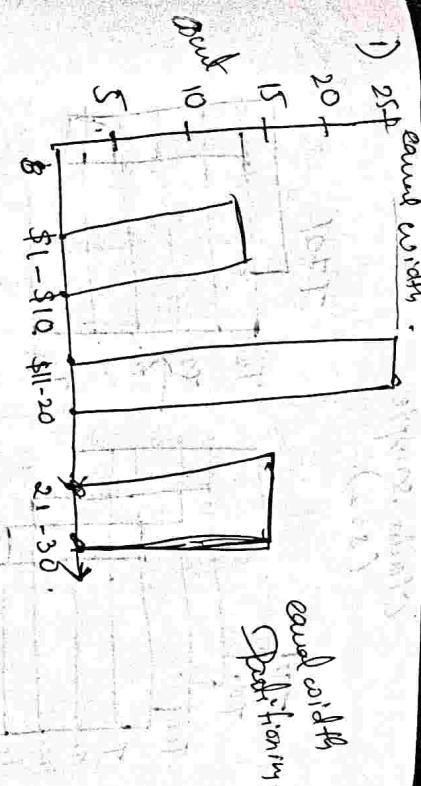
• Simple Random Sampling

• Cluster sample

• Stratified

• Systematic

• Multistage

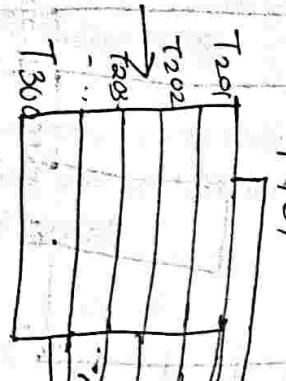
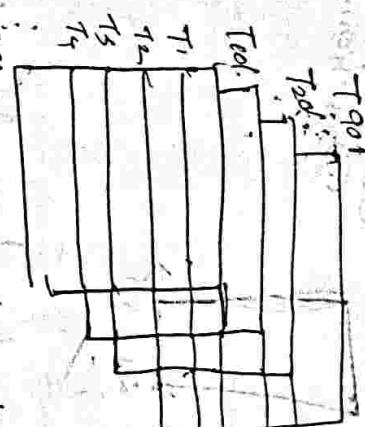


Class Summary

Q3

T701

13/02/2023  
Monday

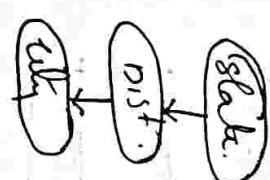


Discretization & concept hierarchy generation

Discretization

Converting continuous data to discrete values.

concept hierarchy.



Note:  
Bar  
bar

bar

bar

2. 1. 11

## Module - 3.

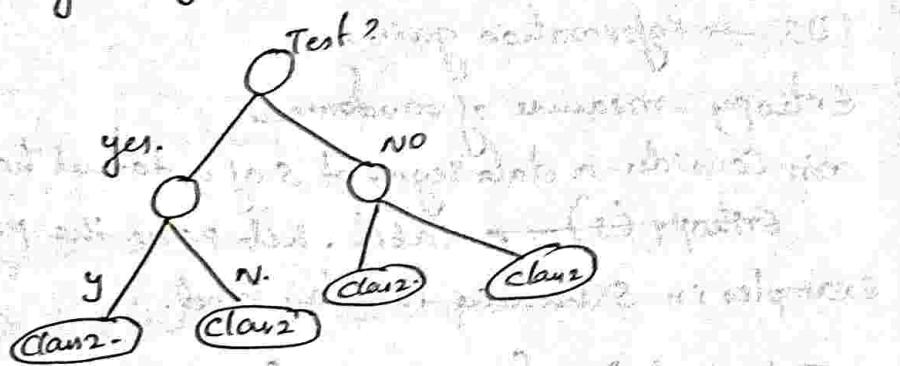
16/12/23 Monday Decision tree

supervised learning method.

It is a hierarchy representation of data.

Simplest type of classification method.

clustering: unsupervised



Nodes: Test on attributes.

Branches: outcome for the test.

Leaves - final classes.

examples of decision tree algorithms.

ID3, C4.5, CART

Two types of decision tree:

1) Classification → o/p are classes.

2) Regression tree (for predicting new value) - leaf nodes are numbers.  
consider the following data set

1)	Name	Features				class label.
		give birth	aquatic animal	aerial animal	has legs	
	human	yes	no	yes	yes	mammal
	Python	no	no	no	no	reptile
	Salmon	no	yes	no	no	fish
	Dog	yes	no	yes	yes	amphibian
	Bat	yes	no	yes	yes	bird
	Pigeon	no	no	yes	yes	bird
	Cat	yes	no	no	no	mammal
	Shark	yes	yes	yes	yes	fish
	Turtle	no	semi	no	yes	amphibian
	Salmander	no	semi	yes	yes	-



$$= -\frac{1}{6} \log_2 \left(\frac{1}{6}\right) - \frac{1}{6} \log_2 \left(\frac{1}{6}\right) - \frac{1}{6} \log_2 \left(\frac{1}{6}\right)$$

$$\frac{3}{6} \log_2 \left(\frac{3}{6}\right) = 1.492$$

$$= \underline{\underline{1.49}}.$$

$$\text{Information gain} = 2.24 - 0 - \underline{\underline{\frac{5}{10}}} \times 1.5 -$$

0.

$$\text{Information gain} = 2.24 - \underline{\underline{\frac{4}{10}}} \times 1.5 - \left(\frac{6}{10}\right) \times 1.492$$

$$= \underline{\underline{0.5712}}$$

$$\text{A} = \text{aerial}$$

$$\text{Value of A} = \{Y, N\}$$

$$\text{Information gain} = \text{Entropy}(S) - \left[ \frac{|S_Y|}{S} \text{Entropy}(S_Y) + \frac{|S_N|}{S} \text{Entropy}(S_N) \right]$$

$$\text{Entropy}(S_N) = \frac{|S_N|}{S} \text{Entropy}(S_N)$$

$$= \frac{3}{6} \text{Entropy}(S_N)$$

$$= 1.905$$

$$\text{Entropy}(S, "aerial animal = no") = \frac{2}{6} \log_2 \left(\frac{2}{6}\right) + \frac{3}{6} \log_2 \left(\frac{3}{6}\right)$$

$$= \underline{\underline{1.49}}$$

$$= \underline{\underline{0}} \quad (\text{pure class})$$

$$\text{Information gain} = \text{Entropy}(S) - \left[ \frac{|S_Y|}{|S|} \text{Entropy}(S_Y) + \frac{|S_N|}{|S|} \text{Entropy}(S_N) \right]$$

$$\text{Entropy}(S, "aerobic animal = no") = P_{\text{mamm}} \log_2 (P_{\text{mamm}}) + P_{\text{bird}} \log_2$$

$$P_{\text{mamm}} - P_{\text{bird}} \log_2 (P_{\text{bird}})$$

$$= -\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) -$$

$$= \underline{\underline{2 \log_2 \left(\frac{2}{3}\right)}}.$$

$$A = \text{has legs.}$$

$$\text{Value of A} = \{Y, N\}$$

$$\text{Entropy}(S, "has legs = yes") = -\frac{2}{7} \log_2 \left(\frac{2}{7}\right) - \frac{3}{7} \log_2$$

$$= \underline{\underline{-\frac{2}{7} \log_2 \left(\frac{2}{7}\right)}}.$$

$$\text{Entropy}(S, "has legs = semi") = \underline{\underline{-\frac{3}{7} \log_2 \left(\frac{3}{7}\right)}}$$

$$= \underline{\underline{-1.5}}.$$

$$\text{Entropy}(S, \text{"has legs = y"}) = -p_{\text{up}} \log_2(p_{\text{up}}) - p_{\text{down}} \log_2(p_{\text{down}})$$

(P.y, P.n)

$$= \frac{1}{3} \log_2 \left( \frac{1}{3} \right) - \frac{2}{3} \log_2 \left( \frac{2}{3} \right)$$

$$= 0.918$$

$$\text{Total Entropy}(S) = \sum -p_i \log_2 p_i$$

$$= -p_{\text{down}} \log_2(p_{\text{down}}) - p_{\text{up}} \log_2(p_{\text{up}})$$

$$= -\frac{5}{10} \log_2 \left( \frac{5}{10} \right) - \frac{5}{10} \log_2 \left( \frac{5}{10} \right)$$

$$= \frac{1}{2}$$

$$\text{Information gain} = \text{Entropy}(S) - \frac{1}{|S'|} \text{Entropy}_{|S'|}$$

$$= 0.24 - \frac{1}{10} \times 1.5 - \frac{3}{10} \times 0.918$$

$$= 0.918 - 0.8856$$

Aquatic animal is the root node because it informs

gain of aquatic animal is higher than non-aquatic animal.

And the root node of Dennis tree may also be animal. Thus, for the following data

Age	competition	Type	claw (proj't.)
old	yes	Software	Down.
old	no	Hardware	Down.
old	no	Hardware	Down.
mid	yes	Software	Up.
mid	no	Hardware	Down.
mid	no	Hardware	Down.
new	yes	Software	Up.
new	no	Software	Up.
new	no	Hardware	Up.

$A = \text{down} \cup \text{competition} \cup \text{age} \cup \text{type}$ .

values of  $A = \{\text{old}, \text{mid}, \text{new}\}$ .

$\text{Entropy}(S, \text{"age = old"}) = -0.0$ .

$$\text{Entropy}(S, \text{"age = mid"}) = -2 \log_2 \left( \frac{2}{4} \right) - 2 \log_2 \left( \frac{2}{4} \right)$$

$$\text{Entropy}(S, \text{"age = new"}) = -1$$

$$\text{Entropy}(S, \text{"age = new"}) = 0$$

$$\text{Information gain} = 1 - 0 - \frac{3}{10} \times 0 = 0.6$$

= 1.

$A = \text{competition}$ .

values of  $A = \{\text{yes}, \text{no}\}$ .

$$\text{Entropy}(S, \text{"competition = yes"}) = -\frac{3}{4} \log_2 \left( \frac{3}{4} \right) - \frac{1}{4} \log_2 \left( \frac{1}{4} \right)$$

$$= 0.8112$$

$$\text{Entropy}(S, "Comp = NO") = -\frac{2}{6} \log_2\left(\frac{2}{6}\right) - \frac{4}{6} \log_2\left(\frac{4}{6}\right)$$

$$= 0.918$$

$$\text{Information gain}_{("Comp")} = 1 - \frac{4}{10} \times 0.8112 - \frac{6}{10} \times 0.918$$

$$= 0.124$$

$A = \text{Software Type}$ .

$$\text{Value of } A = \{S/w, H/w\}$$

$$\text{Entropy}(S, "Type = S/w") = \frac{3}{6} \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \log_2\left(\frac{3}{6}\right)$$

$$= \frac{1}{2} \text{ (prob)}$$

$$\text{Entropy}(S, "Type = H/w") = \frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right)$$

$$(S, \text{Type}) = \frac{1}{4}$$

$$\text{Imp. gain} = 1 - \frac{6}{10} \times 1 - \frac{4}{10} \times 1$$

$$= 1 - 0.6 - 0.4 = 0$$

$$\underline{\underline{.}}$$

Age is the root node.

$$\text{entropy}(A) = H$$

$$A = x_1$$

$$\text{Value of } A = \{T, F\}$$

$$\text{Entropy}(S, "x_1 = T") = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right)$$

$$= 0.918$$

$$\text{Entropy}(S, "x_1 = F") = -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right)$$

$$= 0.916$$

$$\text{Information gain}_{(S, x_1)} = 1 - \frac{3}{6} \times 0.918 - \frac{3}{6} \times 0.916$$

$$= 0.052$$

$$A = x_2$$

$$\text{Value of } A = \{T, F\}$$

Induction no.	Class label	$x_1$	$x_2$
1	T	T	T
2	T	F	T
3	F	T	F
4	F	F	T
5	T	F	F
6	T	T	F

$$= -\frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right)$$

$$= 1$$

$$= -\frac{1}{4} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \times \log_2 \left(\frac{1}{2}\right).$$

$$= 1$$

~~$$\text{Information}(S, \text{age}) = 1 - \frac{3}{4} \times \frac{1}{2} - \frac{1}{2} \times \frac{1}{2}$$~~

$$\text{Gini}(S) = 1 - \sum_{i=1}^C (P_i)^2$$

$$\text{Pamp} = \frac{3}{10}$$

$$\text{Papp} = \frac{1}{10}$$

$$\text{Pbird} = \frac{2}{10}$$

$$\text{Pfish} = \frac{2}{10}$$

$$\text{Pman} = \frac{2}{10}$$

$$\text{Gini}(S) = 1 - \left[ \left( \frac{3}{10} \right)^2 + \left( \frac{1}{10} \right)^2 + \left( \frac{2}{10} \right)^2 + \left( \frac{2}{10} \right)^2 + \left( \frac{2}{10} \right)^2 \right]$$

$$(S) = \frac{0.79}{100}$$

$$\text{Gini split index}$$

Let  $S$  be a set of examples,  $A$  - attribute / feature.

$S_V$  - suitable with  $A = V$

Values(A) - set of all possible values of  $A$ .

Gini split index of an attribute  $A$  is defined by

13/23 Gini Index  
It's a probability of sending a particular class being wrongly classified.

Gini index is used in CART algorithm. Consider a dataset  $S$  having class labels  $C_1, C_2, \dots, C_r$ . Let  $P_j$  is the probability of examples having class label  $C_j$ . Then Gini index is calculated by

$$\text{Gini split}(S, A) = \sum_{\text{values } h} \frac{|S_h|}{|S|} \times \text{gini}(S_h)$$

Gebr. A = ovaatje annies  
Welkes oʃ A = 5 "4en" 41

*Microspilt (S,  $\beta$ -D-<sup>14</sup>C-glucuronate + "green broth").*

A = gives birth  
values of A = "yes", "no".

$$\text{values of } A = \{yes, no\}$$

$$\text{and split}(S, \text{Age} \geq b) = \underbrace{[yes]}_{\text{if } \text{Age} \geq b} \times \underbrace{\text{no}}_{\text{if } \text{Age} < b}$$

$$= \frac{|S_0|}{|S|} \times \min(S_{\text{gambit}})$$

$$\text{Min}(S_{\text{good}} \mid \text{birth} = \text{yes}) = 1 - \sum_i p_i^{1/2}. \quad \text{Total no.}$$

$$P(\text{bind}) = \frac{2}{4}$$

$$P(j|h) = \frac{1}{4}$$

$$= 1 - \left[ \left( \frac{3}{4} \right)^2 + \left( \frac{1}{4} \right)^2 + \left( \frac{1}{4} \right)^2 \right]$$

$$\text{Ans: } (\text{Signs b}^{\circ}\text{ath} = n_0) = 1 - \left[ \left( \frac{1}{6} \right)^2 + \left( \frac{1}{6} \right)^2 + \left( \frac{3}{7} \right)^2 + \left( \frac{1}{7} \right)^2 \right]$$

$$= \overline{0.667} \text{ ft}^3 \text{ min}^{-1}$$

$$\text{CMB split (S, gmc birth)} = \frac{4}{10} \times 6.625 + \frac{6}{10} \times 0.667$$

110.6503

Growth ( $S$ , aquatic animal = year) =  $\frac{d}{2}$

$$\text{min}^i(S, \text{aucutie} = \text{no}) = \text{Eq} \left[ 1 - \left( \frac{\alpha_1}{\beta_1} \right)^2 + \left( \frac{\alpha_2}{\beta_2} \right)^2 + \left( \frac{\alpha_3}{\beta_3} \right)^2 \right]$$

0.64

$$\text{mini\_split(S, max\_cand)} = \frac{\text{size}}{S} \times \text{mini}'(S, \dots)$$

$$+ \frac{[S_{\text{NO}}]}{L_{\text{SI}}} x \sin(S_{\text{NO}} - S_{\text{A}})$$

$$\frac{0.69 \times 5}{10} = 0.32$$

A = aerial animal

values of A = { "yes", "no" }

$$\text{Ovini } (S, \text{serial} = \text{yes}) = 1 - \left[ 1 - \left( \frac{1}{2} \right)^{\frac{1}{n}} \right]^n$$

$$= 0.71875$$

Chini split (S, anisohedral) =

$$\frac{S_{no}}{\text{Sno}} \times \sin^2(\theta_{\text{quad-no}})$$

$$0 + \frac{8}{10} \times 0.71845 \\ = 0.575$$

$A = \text{has legs}$   
Values of  $A = \{\text{"yes", "no"}\}$ .

$$\text{cini}(S, \text{has legs} = \text{yes}) = 1 - \left[ \left( \frac{2}{4} \right)^2 + \left( \frac{2}{4} \right)^2 \right]$$

$$= 0.653$$

$$\text{cini}(S, \text{has legs} = \text{no}) = 1 - \left[ \left( \frac{1}{3} \right)^2 + \left( \frac{2}{3} \right)^2 \right] \\ = 0.44$$

$$\text{cini}_{\text{split}}(S, \text{has legs}) = \frac{|S_{\text{yes}}|}{|S|} \times \text{cini}(S, \text{has legs} = \text{yes}) +$$

$$\frac{|S_{\text{no}}|}{|S|} \times \text{cini}(S, \text{has legs} = \text{no})$$

$$= \frac{7}{10} \times 0.653 + \frac{3}{10} \times 0.44.$$

$$= 0.589$$

split attribute.  
Root node (split node) is selected as the attribute with lowest cini<sub>split</sub> index. Here, animal is split into two split nodes.

One split node is used in each CART algorithm.

$$P_1 = \frac{3}{6} \\ P_2 = \frac{3}{6}$$

$$\text{cini}(S) = 1 - \left[ \left( \frac{2}{6} \right)^2 + \left( \frac{3}{6} \right)^2 \right] \\ = 0.5$$

$A = x_1$   
Values of  $A = \{\text{T, F}\}$ .

$$\text{cini}(S, x_1 = \text{T}) = 1 - \left[ \left( \frac{2}{3} \right)^2 + \left( \frac{1}{3} \right)^2 \right] \\ = 0.44$$

$$= \frac{0.44}{2}$$

$$\text{cini}(S, x_1 = \text{F}) = 1 - \left[ \left( \frac{1}{3} \right)^2 + \left( \frac{2}{3} \right)^2 \right] \\ = \frac{0.44}{2}$$

$$\text{cini}_{\text{split}}(S, x_1) = \frac{|S_{\text{T}}|}{|S|} \times \text{cini}(S, x_1 = \text{T}) + \frac{|S_{\text{F}}|}{|S|} \times \text{cini}(S, x_1 = \text{F})$$

$$= \frac{3}{6} \times 0.44 + \frac{3}{6} \times 0.44$$

$$= 0.44$$

$$= 0.44$$

$A = x_2$   
Values of  $A = \{\text{T, F}\}$ .

$$\text{cini}(S, x_2 = \text{T}) = 1 - \left[ \left( \frac{2}{4} \right)^2 + \left( \frac{2}{4} \right)^2 \right]$$

$$= \frac{0.5}{2}$$

$$\text{cini}(S, x_2 = \text{F}) = 1 - \left[ \left( \frac{1}{2} \right)^2 + \left( \frac{1}{2} \right)^2 \right]. \\ = 0.5$$

$$\text{Gini}(S, X_2) = \frac{4}{6} \times 0.5 + \frac{2}{6} \times 0.0$$

$$= \underline{\underline{0.5}}$$

$\sqrt{0.5}$

$\sqrt{0.5}$  is the root node.

Meaning

### 14.3.3 Performance of a classifier Algorithm

### Tuesday

#### ○ confusion matrix

Actual class	Predicted class	
	Yes	No
Yes	TP	FP
No	FN	TN

2. classification

Probability

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

True positive

False positive

False negative

True negative

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

FP

FN

TN

Actual class

Predicted class

Yes

No

TP

- 3) A database contains 80 reports on a particular type of job which are relevant to certain investigations. A search was conducted on that topic and 50 records were retrieved. 40 were relevant. Construct the confusion matrix for the search and calculate the precision, recall and finally the accuracy whole classification problem

Actual	
relevant	not relevant
predicted relevant	40
not relevant	15
	15

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{TP} = 40$$

$$\text{FP} = 15$$

$$\text{TN} = 15$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{40}{40 + 15} = \frac{40}{55}$$

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}} = \frac{40 + 15}{40 + 15 + 10 + 15} = \frac{55}{80}$$

15/3/23 Wednesday The ID3 Algorithm

Developed by Ross Quinlan  
Alternative Dichotomiser 3

- Assumptions
  - The algorithm uses information gain to select the

most useful attribute for classification.  
Assume that there are only two class labels namely, "+" and "-". The example with class labels "+" are called positive examples and others negative examples.

Notations

The set of examples.

The set of class labels.

The set of features.

An arbitrary feature (attribute).

The set of values of the feature A.

An arbitrary value of A.

The set of examples with  $A = v$ .

Algorithm ID3( $S, F, C$ )

1. Create a root node for the tree.

2. If (all examples in S are positive) then

3. Set the single node tree root with label "+".

4. end if

5. If (all examples are negative) then

6. Set the single node tree root with label "-".

7. end if

8. If (number of feature in S) then

9. Set the single node as a root tree.

10. Else

11. Let A be the feature in F with the highest information gain.

12. Assign A to the root node in decision tree

13. for all  $C$  values,  $v$  of  $N$  do.

14. Add a new tree branch below root node

corresponding to  $v$ .

15. if  $(S_v, \hat{v})$  empty) then .

16. Below this branch add a leaf node with label equal to the most common class label in the set.

17. else

18. Below this branch add the subtree formed by applying the same algorithm 1D3 with the values.

1D3( $S_v, C_v, \hat{v}, \{A\}$ )

19. end if

20. end for

21. end if.

x) consider a two class classification problem of predicting whether a photograph contains a man or a woman. Suppose we have a test dataset of 10 records with expected outcomes and a set of predictions from our classification algorithm.

Expected	Predicted
man	woman
man	man
woman	woman
man	man
woman	woman
woman	woman
woman	woman
man	man
man	woman
woman	woman

10

- a) compute the confusion matrix for the data.  
 b) compute the accuracy, precision recall.

Actual: man	Man	Woman	Actual: woman
	3	1	4
	2	4	5

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{3}{3+1} = 3/4$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{3}{3+2} = 3/5$$

$$\text{Accuracy} = \frac{3+4}{3+2+1+4} = 7/10$$

5) Suppose 10000 patients tested for the disease, 9000 are actually healthy and 1000 are actually sick. For the sick people a test was positive for 620 and negative for 380. For the healthy people, the same test was positive for 180 and negative for 8820. Construct a confusion matrix for the data and compute the accuracy, precision and recall for the data.

Sick	Actual: healthy
620	180
380	8820

10000

$$\text{Predicted} = \frac{620}{620+380} = \frac{620}{1000} = \frac{31}{50}$$

a)  $13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 25, 25, 30, 33, 35, 35, 35, 35, 35, 36, 40, 45, 46, 52, 52$   
 $\text{Mean} = \underline{\underline{29.96}}$

$$\text{Recall} = \frac{620}{620+380} = \frac{620}{1000} = \frac{31}{50}.$$

$$\text{Accuracy} = \frac{620+8820}{620+380+8820} = \frac{9440}{10000} = \underline{\underline{94\%}}$$

b) median = 25.

b) smoothing by bin means

Bin 3: 14, 22, 25, 25.

Bin 1: 13, 15, 16

Bin 2: 16, 19, 20

Bin 3: 20, 21, 22

" 4 : 22, 25, 25.

" 5 : 25, 25, 30.

" 6 : 33, 33, 35

" 7 : 35, 35, 35

" 8 : 36, 40, 45

" 9 : 46, 52, 40.

smoothing by bin mean replaced by median

Bin 1: 14, 6, 14, 6, 14, 6.

Bin 2: 18, 3, 18, 3, 18, 3.

Bin 3: 21, 21, 21

" 4: 24, 24, 24.

" 5: 26, 6, 26, 6, 26, 6.

" 6: 33, 6, 33, 6, 33, 6.

" 7: 35, 35, 35

" 8: 40, 3, 40, 3, 40, 3.

" 9: 56, 56, 56

- c) eliminating outliers (values fall outside the data)
- d) smoothing by bin median  
 " " boundary
- e) use min-max algorithm normalization to transform value 35 for age from 0 to 100  
 $[0:100, 1:0]$
- f) use Z-score normalization to transform the value 35 for age, where the standard deviation of age is 12.94 (mean 35)  
 g) use normalization by decimal scaling to transform the value 35
- h) convert or which method of normalization you would prefer to use for the given data? giving reasons as to why?
- i) use Z-score normalization to transform the value 35 for age, where the standard deviation of age is 12.94 (mean 35)  
 j) use normalization by decimal scaling to transform the value 35

(e)  $\text{new\_max} = 1.0$

$$n_{\text{res}, \text{max}} = 1.0$$

$$\sqrt{35} = 13$$

$$V^* = \min_{\pi} V^\pi$$

MAXIMA-MINIMA

$$= \frac{35 - 13}{(1.0 - 0)} + 0.$$

二〇三

11

$\sigma_A$

$$\bar{A} = 29.96.$$

$$\frac{V' = 35 - 29.96}{12.94}$$

$$= \underline{0.389}$$

$$g) \sqrt{1} = \frac{\sqrt{35}}{10} = \frac{35}{100} = \underline{\underline{0.35}}$$

h) Z-score normalization becomes a good method

Suppose a hospital treated the 200,000,000 Indians  
of India for 10 years only selected adults, with one doctor

5/8  
32.9 41.2 35.4  
1. calculate the mean, median and SD of age and fat  
2) calculate the two variables based on z-score normalization  
3) normalize the two variables (Pearson's product-moment correlation coefficient)  
4) calculate the correlation coefficient (positive or negative)

$$\text{Age} = \frac{46+4}{2} = 45$$

1.8, 9.5, 14.8, 25.9, 26.5, 27.

$$\text{Mean of first} = \frac{26.9}{1}$$