

S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	S	M	T	W	T	F	S	S	M					
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30

117-248  
18th WeekData:-

Thursday

27

It is a raw fact, information or statistics which can be in various forms like numbers, texts or any other form.

Why Data Mining?

→ The explosive growth of data means production of data is too much in various areas. Due to increase in size of database, increase in computerized growth in the society, automatic analysis overcome on manual analysis which shows need of data mining.

1960 → Data collection and Database creation

1970 → Database Management System (SQL)

1980 → ~~RDBMS~~ RDBMS i.e. advancement of DBMS

1990 - 2000 → Data Mining

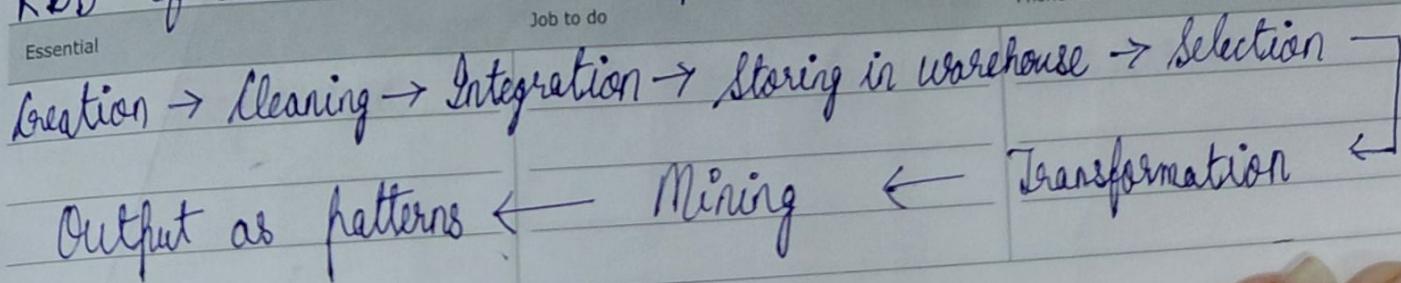
What is Data Mining?

→ It refers to extracting or mining knowledge from large datasets  
 → It is also known as knowledge discovery <sup>from</sup> ~~of~~ data.

KDD of data involves 7 steps:-

Job to do

Phone No.



S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M					
1 11-246 18th Week	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17 21 24	18 22 25	19 23 26	20 27 28	21 29	30

Data Cleaning :- To remove the Noise and inconsistent data Saturday

29

Data Integration :- Where multiple data sources may be combined and stored in data warehouse.

Data Selection :- Where data relevant to the analysis are retrieved from database.

Data Transformation :- Where data are transformed into appropriate forms for mining by performing scaling operations.

Data Mining :- An essential process where intelligent methods are applied in order to extract data patterns.

Pattern Evaluation :- To identify interesting patterns, representing knowledge based on measures

30 Sunday

Knowledge Presentation :- Where visualisation and knowledge representation techniques are used.

Eg:- Pie chart, Bar chart, Area under curve, etc.

## Architecture of Data Mining

i) Data collected goes through cleaning, integration and selection. (Preprocessing)

- |           |           |           |
|-----------|-----------|-----------|
| Essential | Job to do | Phone No. |
|           |           |           |
|           |           |           |
|           |           |           |
|           |           |           |
|           |           |           |
- ii) The preprocessed data is stored in warehouse server.

O | (iii) The intelligent algorithms are applied through data mining.

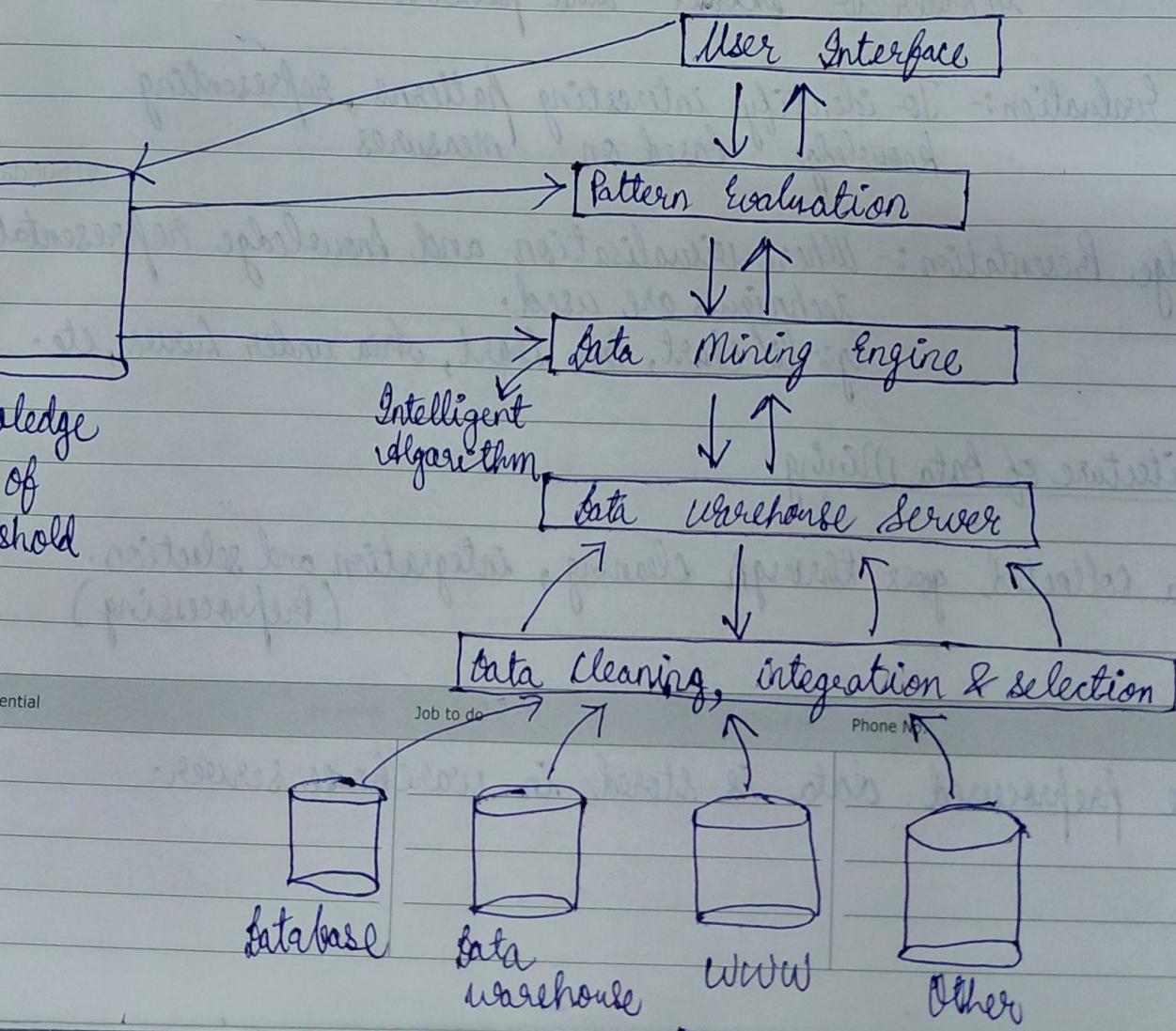
121-244  
19th Week

(iv) The mined patterns are evaluated, output is measured.

For finding out whether the output is accepted by user we require background knowledge!

The threshold in the extra knowledge, if greater than output, the output is acceptable by the user.  
If threshold < output, output is not accepted.

Only the user can change the threshold value.



M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

122-243  
19th Week

→ Database, Data Warehouse, WWW, other repositories

Tuesday

# 02

↳ These are set of database, datawarehouse, spreadsheet or other kind of information repositories.

→ The Database or Data warehouse server is responsible for fetching the relevant data based on user requests.

→ Knowledge Base

↳ This is the background knowledge used to guide the search engine or evaluate the resulting patterns.

→ Data Mining Engine

↳ This consists of a set of modules for task prediction and evaluation analysis.

→ Pattern Evaluation

↳ This component interact with data mining modules based on the search.

→ User Interface

↳ This module communicates between user and the data mining system. It allows user to interact with system by specifying data mining query.

Essential

Job to do

Phone No.

03

Wednesday

Predictive analysis123-242  
19th Week

- It is the process of extracting information from large dataset in order to make prediction and estimates about future outcomes.

Eg:- Loan Default Predictions

Data Mining Knowledge - Analysing past records, the bank finds customers with high credit utilization and irregular payment history are likely to default.

When a new customer applies for a loan the bank uses this pattern to credit, the likelihood of default.

Eg:- Health Care - predicting a patient prone to heart disease.

What kind of data?

- Number of data repositories on which mining can be performed.

Repositories are:-

↳ Relational Database

↳ Data Warehouse

↳ Transactional Database

↳ Advanced Database System

↳ World Wide Web

→ Object Oriented / Object Relational Database

→ Spatial Database

→ Temporal Database, Sequential Database

→ Text Database, Multimedia Database

→ Heterogeneous Database and Legacy Database

→ Time Series Database

M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

124-241  
19th Week

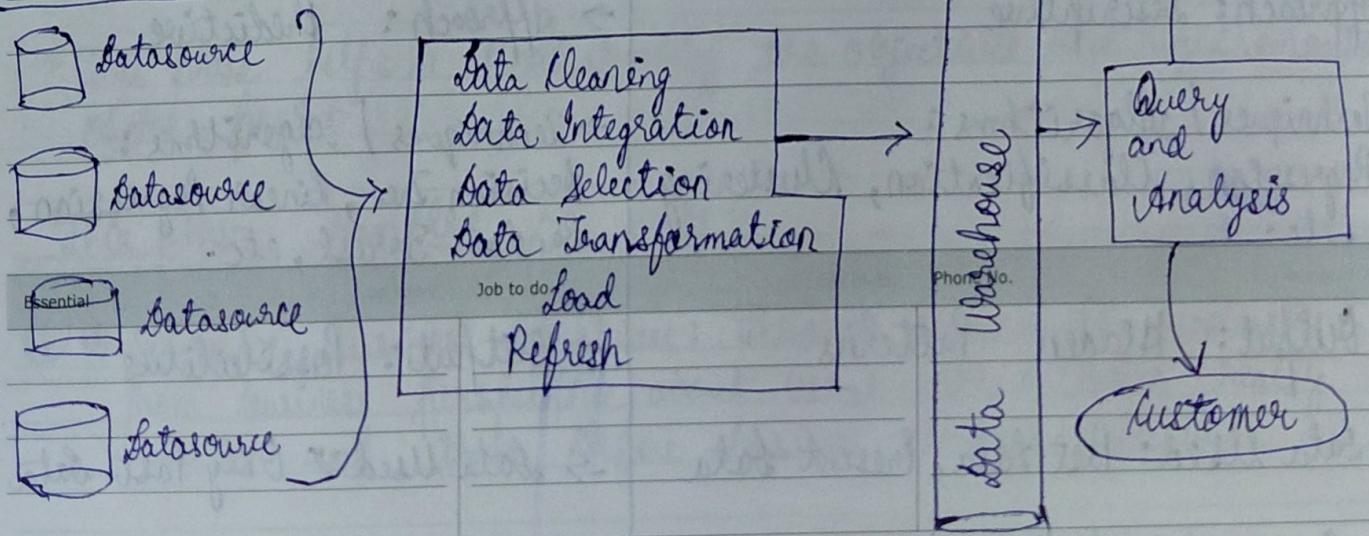
## Relational Database

Thursday

04

- DBMS consists of database and set of software programs
- The database consists of a collection of interrelated data
- Software programs are used to manage and access the data, stored in database
- Software programs consists of database structure, data storing, distributing data access and for ensuring consistency and security of the information stored.
- The Relational Database is a collection of tables, each table is assigned an unit name.
- Each table, consists of set of attributes and large set of tuples.
- Each tuple is represented by a unique key and described by a set of attributes.

## Data Warehouse



05

→ A data warehouse is a ~~self~~ repository of information collected from multiple sources ~~stored~~ stored under unified scheme and usually resides at a single point.

Friday

125-240  
19th Week

- Data warehouse are constructed through a process of data cleaning, data integration, data transformation, data loading and periodic data ~~refreshment~~ refreshing.
- In computing, data warehouse is also known as enterprise data warehouse (EDW). It is a system used for reporting and data analysis.

### \* Differentiate between Data Mining and Predictive Analysis

#### Data Mining

- It is a process of extracting ~~know~~ knowledge from large database
- Objective: To explore and understand data
- Approach: Descriptive
- Techniques / Algorithms: Regression, Classification, Clustering, etc.
- Output: Hidden Patterns
- Data Used: Past data, Present data
- It is independent

#### Predictive Analysis

- Applying the data mining knowledge or patterns to predict future outcome.
- Objective: To predict future values and possibilities
- Approach: Predictive
- Techniques / Algorithms: Decision Tree, Linear Regression, Random Forest, etc.
- Output: Possibilities
- Data Used: Only Past data
- It is dependent

M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

126-239  
19th Week

## CRISP - DM :-

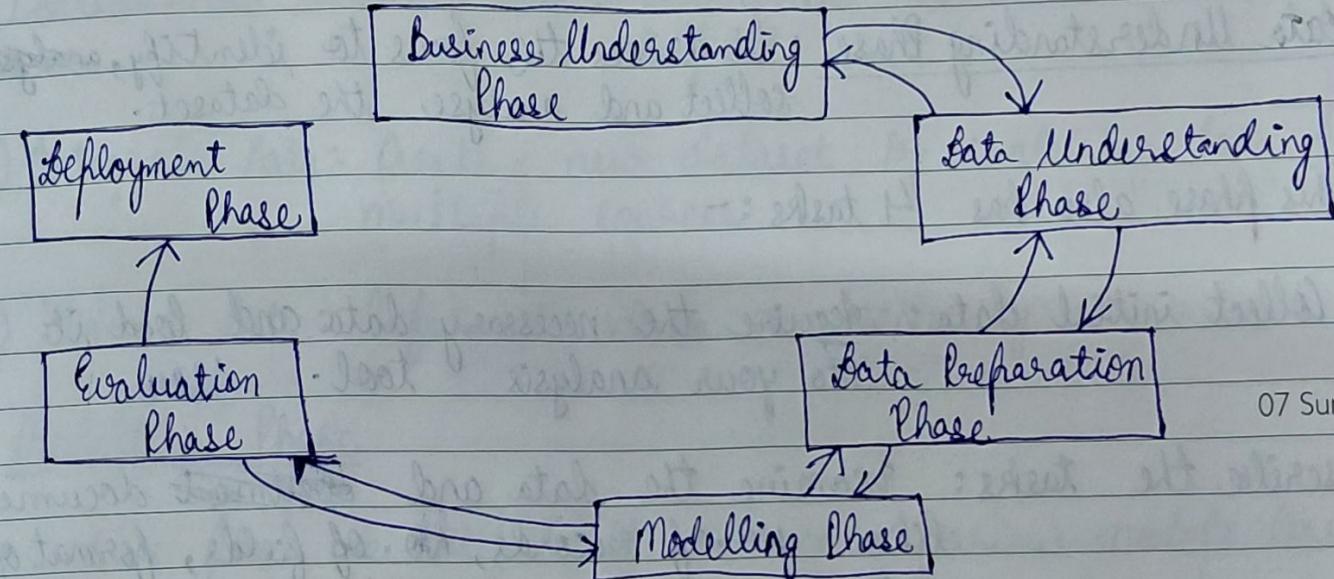
Saturday

06

"Cross Industry Standard Process for Data Mining"

CRISP provides standardised data mining process across industries for fitting data mining into general problem solving strategy of a business or research unit.

It has 6 phases :-



07 Sunday

### Business Understanding Phase

→ This phase helps in understanding the objectives and requirements of the project.

This phase handles 4 tasks :-

- | Essential | Job to do   | Phone No. |
|-----------|---|-----------|
| ①         | Determine Business Objectives : Through thoroughly understand from business perspective about what the customer really wants and then define business success criteria. |           |

08

Monday **(ii)** Access Situation: Determine Resources availability, project requirement, access contingency and conduct a cost benefit analysis.

- (iii)** Determine Data Mining goals: Define technical Data Mining perspective
- (iv)** Product Project Plan: Select technologies and tools and define detailed plan for each project phase.

Data Understanding Phase: It drives the focus to identify, collect and analyse the dataset.

This phase also has 4 tasks:-

- i** Collect initial data: Acquire the necessary data and load it into your analysis tool.
- ii** Describe the tasks: Examine the data and ~~document~~ document, like No. of records, No. of fields, format of data, etc.
- iii** Explore data: Visualise and Identify relationship among data
- iv** Verify data: Clear or first data must have to be verified.

Essential

Job to do

Phone No.

M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

129-236  
20th Week

## Data Preparation Phase

Tuesday

09

It prepares the final dataset for modelling

Tasks:-

- i) Select Data: Choose the dataset that will be used
- ii) Clean Data: To correct or remove erroneous data.
- iii) ~~Derive New~~ Construct Data: Derive new attribute that will be helpful
- iv) Integrate Data: Create a new dataset by combining data from multiple sources.
- v) Format Data: Reformat data as necessary

## Modelling Phase

→ In this phase we will build and access various models based on several modelling techniques.

Tasks:-

- i) Select Modelling Technique: determine which algorithm you want to try

Essential	Job to do	Phone No.
ii)	Generate Test Designs: Needs to split the data into training, testing and validation.	

10

Wednesday

130-235  
20th Week

(iii) ~~Create~~ Build Model: Execute few lines of code

(iv) Access Model: Multiple models are competing against each other to interpret model results based on domain knowledge.

### Evaluation Phase

It focuses on technical model assessment.

Tasks:-

i) Evaluate Results: Which one we approve for the business.

ii) Review process: Check the model is properly executed or not, summarise it and correct anything if needed.

iii) Determine next phase: Based on previous tasks, determine whether to ~~deploy~~ proceed for deployment or not.

### Deployment Phase

A model is not particularly useful unless the customer can access ~~the~~ its result.

Essential	Job to do	Phone No.

M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	F	S	P
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

131-234  
20th Week

Tasks:-

Thursday

- i) Plan Deployment : Develop and document a plan for deploying the model
- ii) Plan Monitoring and Maintenance - to avoid issues during the ~~operation~~ phase of model.
- iii) Produce final report = Document a summary of the project
- iv) Review Project - What went well, what would have been better, How to improve

### Fallacies of Data Mining (Misconceptions)

- i) Data Mining tools are automated tools that can be deployed on data repositories to find and search our problems.
- Reality: There are no automatic data Mining tools which will automatically solve your problems.
- ii) Data mining process is autonomous requiring no human intervention
- Reality: Without experts, blind use of data mining software will only provide wrong answers to wrong questions.

Essential	Job to do	Phone No.
iii)	Data Mining pays for itself quickly	
Reality:	Benefits take time. Costs include startup costs, data collection, data warehouse preparation costs, etc	

12

iv Data Mining software packages are easy to use.  
 Friday

132-23:  
20th Week

Reality : No ease of use varies it may require statistical and mathematical knowledge of particular application domain.

v) Data Mining will identify the causes of other business problems or research problems.

Reality : It will help to show hidden patterns but human judgement and domain expertise is needed.

vi) Data Mining will automatically clean up messy database

Reality : Still needs to handle missing values, outliers and inconsistency ~~manually~~ manually

vii) Data Mining always provides positive results

Reality : Not guaranteed. Poor data or assumption can lead to misleading outcomes.


M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

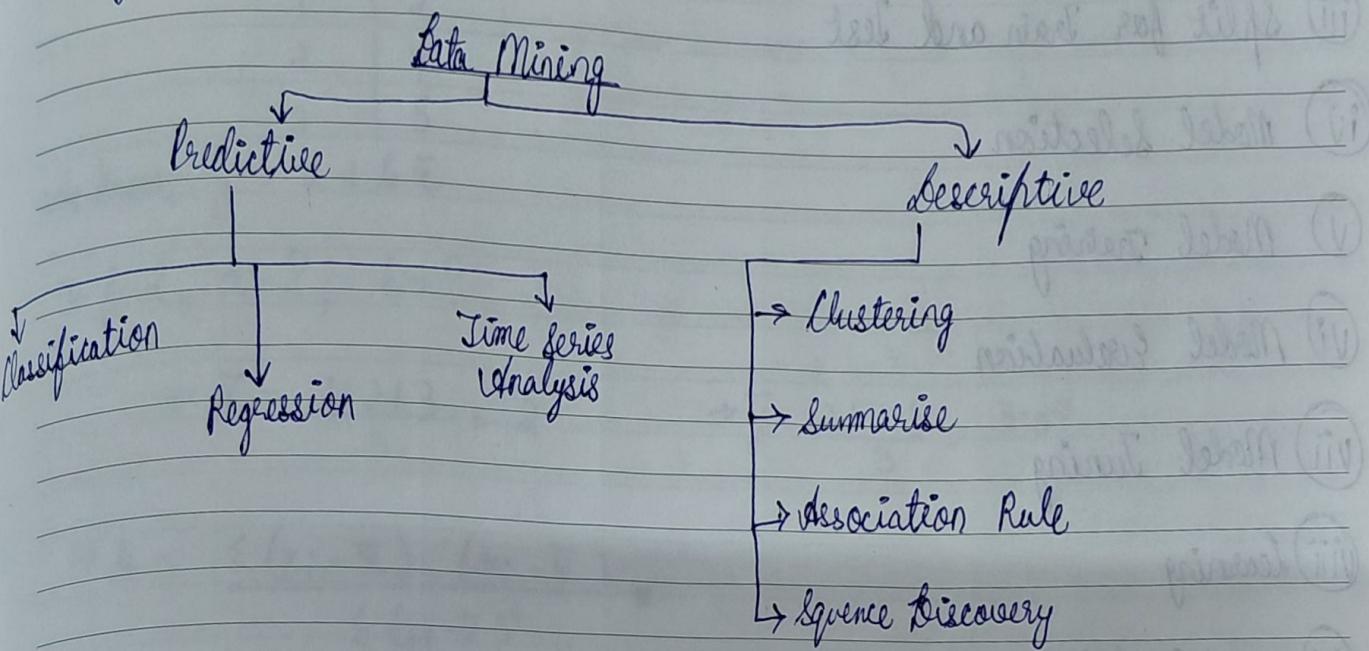
133-232  
20th Week

## Tasks of Data Mining

Saturday

13

Data Mining tasks involve finding patterns and useful information from large dataset.



## Classification

14 Sunday

- It is a ~~pre defined~~ supervised learning technique in data mining that involves categorising or classifying data in pre defined classes.
- Categorising all groups based on features or attributes
- In supervised learning, labelled data can be used to build a model that can predict the class of unseen data.

Essential → It is of two types <sup>to do</sup> Binary Classification and Multiclassification

MAY

2023

15

Steps to build Classification:

Monday

135-230  
21st Week

- (i) Data Preparation
- (ii) Model Feature Selection
- (iii) Split for Train and Test
- (iv) Model Selection
- (v) Model Training
- (vi) Model Evaluation
- (vii) Model Tuning
- (viii) Learning
- (ix) Model Deployment

### Regression (Estimation)

→ It is a supervised learning technique used to predict a continuous numerical value by analysing past data.

Types:-

- (i) Linear Regression    (ii) Logistic Regression    (iii) Polynomial Regression

Essential	Job to do	Phone No.
(iv) Lasso Regression	(v) Ridge Regression	

M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

136-229  
21st Week

Q) We have a dataset Hours studied ( $x$ ) to predict exams( $y$ ) Tuesday  
We need to find  $y$  if  $x=4$

16

$x$	$y$
1	2
2	4
3	5

We know,  $\hat{y} = a + b\bar{x}$

To find,  $a = ?$ ,  $b = ?$

$$\rightarrow \bar{x} = \frac{1+2+3}{3} = 2$$

$$\rightarrow \bar{y} = \frac{2+4+5}{3} = 3.67$$

$$\rightarrow b = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$x$	$y$	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	2	-1	-1.67	1.67	1
2	4	0	1.67	0	0
3	5	1	2.67	2.67	1

$$\therefore b = \frac{1.67 + 2.67}{1+1} = \frac{3.67}{2} \rightarrow b = 1.83$$

$$a = \bar{y} - b\bar{x} = 0.67$$

$$\text{Test Data} \Rightarrow \bar{y} = 0.67 + 1.83x$$

Essential

Job to do

Phone No.

So, for  $x=4$ ,

$y = 6.67$

17

Prediction

Wednesday

137-228  
21st Week

→ This is similar to classification and estimation except for prediction, the result lie in future.

Eg:- Predicting prices of stock 3 months into the future.

Time series analysis

→ It is a way of analysing sequence of data points connected over an interval of time.

It is composed of 4 main elements :-

- i Trend: Long term movement of data
- ii Seasonality: predictable recurring patterns of repetition over a fixed regular interval.
- iii Cycle: Long term Movement or fluctuation
- iv Irregularities: Random or unpredictable movement in data

Clustering: It means grouping of records, observations into a similar object.

This is mainly used in unsupervised learning. The data points

having less distance belongs to same cluster and having more distance belongs to different clusters.

M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

138-227  
21st Week Clustering segments the whole dataset into relatively homogeneous subgroups. Thursday

18

Summarisation : It is the process of reducing large dataset into shorter, more understandable format that highlight key patterns, trends and relationships.

Types:-

↳ Descriptive : Uses statistical measures to describe main features of numerical data.

Eg:- Mean, Median, Mode

↳ Aggregation : Combining data into simpler summarised form like sum, average, etc.

↳ Sampling : Select a representative subset of the data.

Association : Association for data mining is the job of finding which attributes go together.

Eg:- Market basket analysis

This rule shows how frequently an item set occurs in a transaction.

Sequential Discovery : It is a data mining technique used to

Essential

Job to do

Phone No.

identify frequently occurring patterns in sequential data such as purchase history.

MAY

2023

19

Friday

## APPENDIX : Data Summarisation and Visualisation

Applicant	Martial Status	Mortgage	Income	Rank	Year	Risk
1	Single	Y	38K	2	2009	Good
2	Married	Y	32K	7	2010	Good
3	Other	N	25K	3	2011	Good
4	Other	N	36K	3	2009	Good
5	Other	Y	33K	4	2010	Good
6	Other	N	24K	10	2007	Bad
7	Married	Y	25K	8	2009	Good
8	Married	Y	48K	1	2010	Good
9	Married	Y	32K	6	2011	Bad
10	Married	Y	32K	5	2009	Good

Variables

↓  
Qualitative  
(Categorical)

↓  
Quantitative  
(Numerical)

Nominal

Ordinal

Interval

Ratio

Discrete Value

Continuous Value

Essential

Job to do

Phone No.

M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

140-225  
21st Week

20

Qualitative :- Describes qualities or categories of data Saturday  
It is also called categorical variables.

Eg:- Marital status, Mortgage, Rank, Risk.

Quantitative :- Takes numerical values and allows arithmetic operation. It is also called numerical variables.

Eg:- Year, Income.

Nominal :- Used for names, labels or categories without order

Eg:- Marital status, Mortgage, Risk.  
(Blood groups)

Ordinal :- Maintains a particular order Eg:- Rank

Interval :- It is a quantitative data defined on an interval without a natural zero (0).

21 Sunda

Eg:- Year

Ratio :- Quantitative data for which mathematical operations can be performed (can contain 0)

Eg:- Income

Essential

Job to do

Phone No.

22

Monday

discrete Values :- It is a quantitative variable that can take finite or countable numbers including zero.

142-223  
22nd Week

- It cannot take values in between two numbers → Gaps between values
- Graphical Representation is Bar Chart.

Eg :- Year

Continuous Values :- It is a quantitative variable that can take any value within a given range.

- No Gaps between the points. → Eg :- Income
- Graphical Representation - Histogram

Predictor Variable : It is a variable whose value is used to predict the value of response variable

The predictor variable in table are all variable except Risk.

Response Variable : It is also called dependent variable or output variable or target variable which try to predict, explain or measure the effect.

Eg :- Effect of 'x' hours studies on exam score 'Y'. Example :-  
In table Risk

Predictor = X Job to do

Phone No.

Response = Y.

M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

143-222  
22nd Week

## Measurement of center, variability and position

Tuesday

23

→ Measures of Central tendency - It measures the location of middle or center of a data distribution. Eg:- Mean, Median, Mode and Midrange

Mean - It is an arithmetic average. The sum of all data values divided by the number of values.

Eg:- Dataset → 2, 4, 6, 8, 10

$$\text{Mean} = \frac{2+4+6+8+10}{5} = 6$$

Advantage :- Easy to compute including all data points

Disadvantages :- Sensitive to outliers

Median - The middle value when the data is sorted.

Eg:- 3, 5, 7, 9, 11

$$\text{Median} = 7$$

3, 5, 7, 9

$$\text{Median} = \frac{5+7}{2} = 6$$

Advantage: Not affected by outliers

Mode - The value that appears most frequently in the dataset.

It is of 3 types: Unimodal, Bimodal and Multimodal

Essential No Mode - If all values occur with same frequency

Eg:- 2, 3, 3, 4, 4, 5, 5, 8, 8, 8, 10, 12

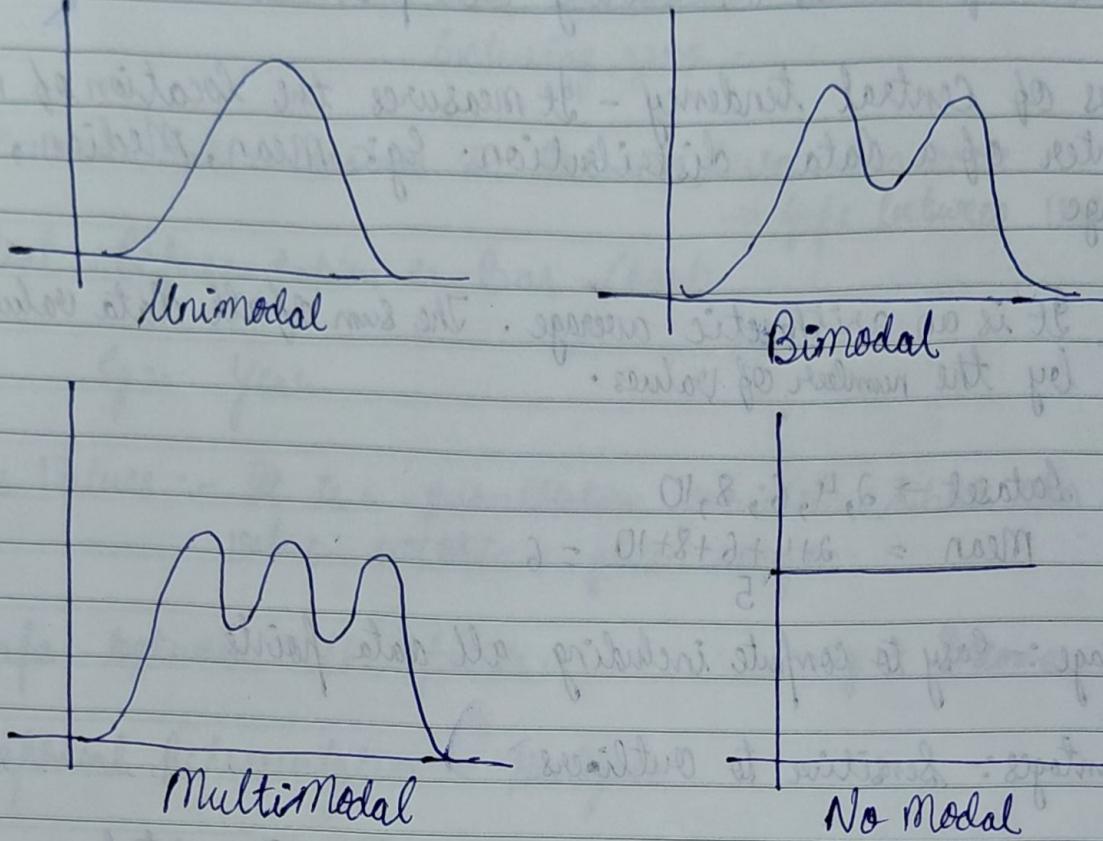
$$\text{Mode} = 8$$

MAY

2023

24

Wednesday

144-221  
22nd Week

**Advantage :-** It is used mostly for categorical data

**Midrange -** It is the average of largest and smallest value in dataset.  
It is applicable for numerical data.

Measure	Best Used	Sensitive to outliers	Categories
Mean	Numerical data and symmetric data	Yes	Interval / Ratio
Median	Skewed data or weak outliers	No	Ordinal and Interval / Ratio
Mode	Most frequent item needed	No	Nominal, ordinal, Interval / <del>Ratio</del> Ratio

M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

145-220  
22nd Week

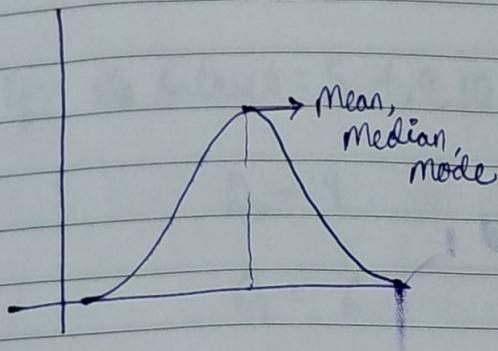
## Measuring dispersion of data

Thursday

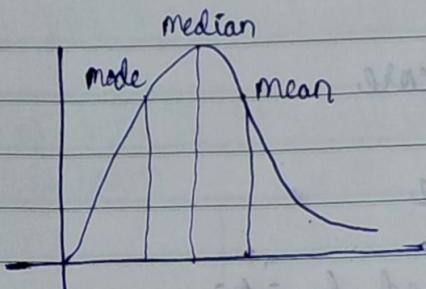
25

In unimodal, frequency curve with perfect symmetric data distribution - mean, median and mode are all at same center value.

But, in reality, data is not symmetric

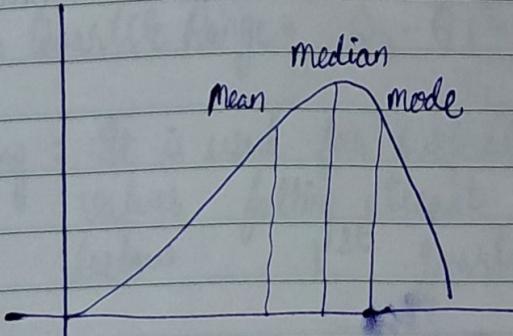


Positive Skewed:



(mode < median)

Negative Skewed:



(Median < Mode)

Essential

Job to do

Phone No.

MAY

2023

26

Friday

## Measurement of variability dispersion (variability)

146-219  
22nd Week

- i Range
- ii Quartiles ( $Q_1, Q_2, Q_3$ )
- iii Inter quartile range (IQR)  $\rightarrow Q_3 - Q_1$
- iv 5 Number Summary ( $Q_1, Q_2, Q_3, \text{min}, \text{max}$ )
- v Box Plot
- vi Z score
- vii Variance
- viii Standard deviation

Range : It is the difference between maximum and minimum value in dataset.

Eg:- Dataset : 5, 7, 9, 10, 12

$$\text{Range} = 12 - 5 = 7$$

Essential

Job to do

Phone No.

M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

147-218  
22nd Week

Quartile :- It divides ordered data into 4 equal parts Saturday

27

$Q_1 \rightarrow$  First quartile, which represent 25% ile of data below it

,  $Q_2 \rightarrow$  Second quartile, known as median, 50% ile of data below it

$Q_3 \rightarrow$  Third quartile, which present 75% ile of data below it

Eg:- Dataset : 5, 7, 9, 10, 12

$$Q_2 = 9$$

$$Q_1 = 6$$

$$Q_3 = 11$$

$$\text{Inter Quartile Range} = Q_3 - Q_1 = 5$$

28 Sunday

Lapping : It is used for suspected outliers, common rule is single out values falling atleast  $1.5 \times IQR$  (lapping) above 3rd quartile or below 1st quartile.

Five Number Summary : It is a concise statistical summary of data consisting of min,  $Q_1$ ,  $Q_2$ ,  $Q_3$ , max

Dataset : 5, 7, 8, 12, 13, 14, 18, 21, 23, 25

Essential

Job to do

Phone No.

$$\text{Min} = 5$$

$$\text{Max} = 25$$

$$Q_2 = 13.5$$

$$Q_1 = 8$$

$$Q_3 = 21$$

29

Monday

## Box Plot (Whisker Plot)

149-216  
23rd Week

→ It is a graphical representation of 5 Number summary of a dataset.

→ Features:

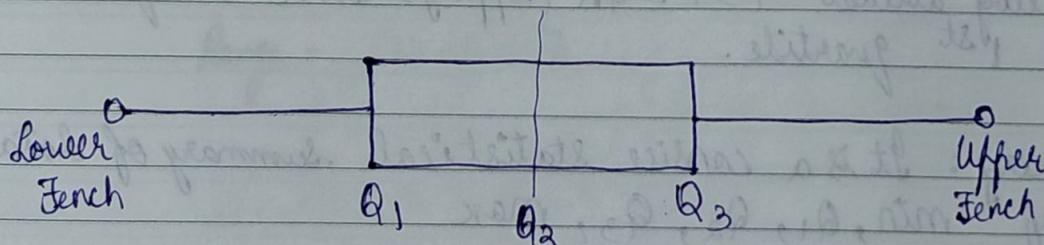
- i) Box: Extends from  $Q_1$  to  $Q_3$
- ii) Median line: A line inside the box shows median
- iii) Whisker: A line extending from  $Q_1$  to minimum and  $Q_3$  to maximum.

Outliers

→ Data points that lie beyond lower fence and upper fence

$$\text{Lower fence} = Q_1 - 1.5 \times \text{IQR}$$

$$\text{Upper fence} = Q_3 + 1.5 \times \text{IQR}$$



Box Plot

M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

150-215  
23rd Week

Q) In dataset 11, 7, 35, 55, 64, 86, 88, 90, 95, 97 Tuesday

30

Sorted = 7, 11, 35, 55, 64, 86, 88, 90, 95, 97

$$\text{min} = 7$$

$$\text{max} = 97$$

$$Q_2 = 75$$

$$Q_1 = 35$$

$$Q_3 = 90$$

$$\text{IQR} = 55$$

$$\text{Capping} = 1.5 \times 55 = 82.5$$

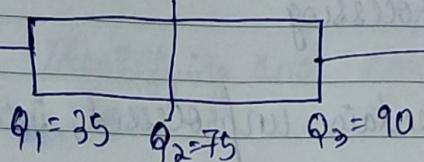
$$\text{Lower} = -47.5$$

$$\text{Upper} = 172.5$$

so, No outliers

$$\text{Min} = 7$$

$$\text{Max} = 97$$



## Z score

→ Z score measures how many standard deviations a data point is from the mean. It standardises data making values from different distribution comparable.

$$Z = \frac{x - \mu}{\sigma}$$

If  $Z = 0$  → Value is exactly at mean  
below

If  $Z < 0$  → Value is ~~above~~ mean

Essential

Job to do

Phone No.

If  $Z > 0$  → Value ~~is~~ is above mean

→ If datapoint with  $|Z| > 3$  are considered outliers

31

Wednesday Variance

- It measures the average square deviation from the mean

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

### Standard Deviation

- It is the square root of variance  $\sqrt{\sigma^2} = SD$

### Chapter-2 Data Preprocessing

- Raw data in database contain unprocessed, incomplete and noisy data.

- Data preprocessing is the process of preparing raw data for analysis by cleaning and transforming it into a usable format.

- Real dataset often have missing or incomplete entries. Data may contain type error, duplicates or outliers.

- Different features may have different scale not suitable for data mining models.

- Some features are irrelevant, redundant or highly correlated.

Essential

Job to do

Phone No.

- Values are not consistent with the policy.

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25

152-213  
23rd Week

## Steps of data processing:-

Thursday

01

- i) Data Cleaning
- ii) Data Integration
- iii) Data Transformation
- iv) Data Reduction

### Data Cleaning

- It is the process of identifying and correcting errors and inconsistencies in the ~~whole~~ dataset.
- It involves handling missing values, removing duplicates and correcting incorrect or outlier data to ensure the dataset is accurate.
- Clean data is ~~is~~ essential for effective analysis as it improves the quality of results and enhance the performance of data models.

### Missing Values

- When a data is absent in a dataset is called missing values.

### Handling of Missing Values

Essential

Job to do

Phone No.

- a) Ignore the tuples or records : This is usually done when class label is missing.

02

Friday This method is not very effective unless the tuple contains several attributes with missing values.

153-212  
23rd Week

b) Fill in the missing value manually

This method is time consuming may not be feasible for large dataset.

c) Use a global constant to fill in the missing values

d) Use the attribute mean to fill in the missing value-

e) Use the attribute mean for all the samples belonging to the same class as the given tuple.

f) Use the most preferable probable value to fill in the missing value.

Eg :- Decision Tree, Regression --

### Noisy Data

→ Noise is a random error or variance in a measured variable.

→ Incorrect attribute values may be due to faulty data, data entry problem, data transmission problem, technology limitations, inconsistency in

Essential

Job to do

Phone No.

Naming convention.

→ Other problems include duplicate record, incomplete data

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25

154.211  
23rd week

## Binning Method

Saturday

03

Sort the data and the sorted values are distributed into a number of buckets or bins. Then these bins can be smoothened up by mean, median and boundaries.

### Bin Mean

→ Each value in bin is replaced by mean value of the bin

### Bin Median

→ Each bin value is replaced by the bin median

### Bin Boundary

→ Maximum and Minimum values in a given bin are identified as bin boundaries.

→ Each bin value is then replaced by the closest boundary value.

### Equal Width Partitioning

→ It divides the range into N intervals of equal size.

→ If A and B are lowest and highest value of the attribute the width of the interval will be  $\frac{B-A}{n} = w$

Essential This is

Job to do

Phone No.

→ The most straight forward method but outliers may dominate presentation.

JUNE

2023

05

Monday

Equal Depth Partitioning (Equal Frequency Partitioning)

→ It divides the range into  $N$  intervals each containing approximately same number of samples.

good data scaling

↳ Managing Categorical attributes

Q) Sorted data for price is 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34

Smoothing By Bin Means

Mean = 9

22.75

29.25

Bin-1 - 9, 9, 9, 9

Bin-2 = 22.75, 22.75,

22.75, 22.75

Bin-3 = 29.25, 29.25

29.25, 29.25

Smoothing By Medians

Bin-1 > 8.5

Bin-2 > 22.5

Bin-3 > 29.25

Smoothing By Boundaries

Bin-1 > 4

Bin-2 > 21

Bin-3 > 34

Essential

Job to do

Phone No.

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S						
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30

157-208  
24th Week

## Equal Width Partitioning

Tuesday

06

~~Step-1~~  $\text{Min} = 4$

$\text{Max} = 34$

$\text{Range} = 34 - 4 = 30$

~~Step-2~~  $\text{Width} = \text{Interval} = \frac{\text{Range}}{\text{Intervals}} = \frac{30}{3} = 10$

No. of Bins  
(default 3)

~~Step-3~~ Bin-1  $[4, 4+10] = [4, 14] = 4, 8, 9$

Bin-2  $[14, 14+10] = [14, 24] = 15, 21, 21, 24$

Bin-3  $[24, 34] = 25, 26, 28, 29, 34$

## Equal Depth / Frequency partitioning

~~Step-1~~ Data must be sorted

~~Step-2~~ Take 3 Bins . Total value in dataset is 12. So, each bin should have 4 values

~~Step-3~~ Bin-1) 4, 8, 9, 15

Bin-2) 21, 21, 24, 25

Bin-3) 26, 28, 29, 34

Essential

Job to do

Phone No.

JUNE

2023

07

Wednesday

## Regression

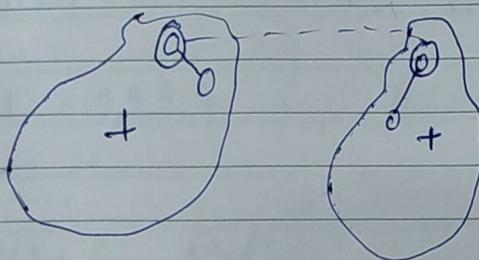
→ In Regression, data can be smooth by fitting the data to a function.

→ Two types :-

i) Linear Regression

ii) Multiple Linear Regression

## Clustering



## Miss Classification

→ It means a model predict the wrong class level compared to the actual or true class level in your dataset.

## Dataset

	<u>Brand</u>	<u>Frequency</u>	Job to do	Phone No.
Europe	US X	USA	1	
	France		1	
	US	156		
	Europe	46		
	Japan	51		

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	.

159-206  
24th Week

# 08

- In frequency distribution dataset, having 5 classes ~~now~~ Thursday  
 however two of the class USA and France have count only one  
 each which is clearly happening that two of the records have been  
 inconsistently classified (Misclassification) with respect to  
 original data.
- To maintain the consistency, record have been labelled US and  
 Europe instead of USA and France.

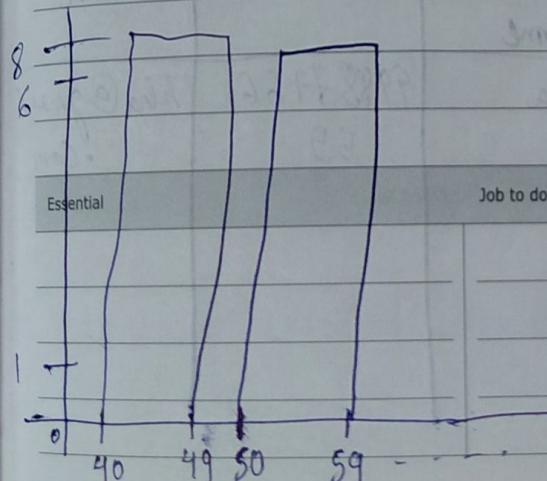
## Graphical Methods for identifying outliers

Outliers are extreme values that go against the trend of remaining data.

Graphical Methods for identifying outliers in numerical variables  
 are Histogram and Scatterplot.

Dataset : 45, 50, 52, 48, 47, 49, 51, 46, 500, 50, 48, 49, 52,  
 47, 51

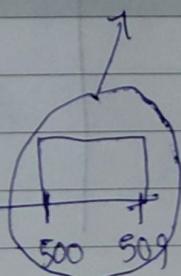
Sorted : 45, 46, 47, 48, 49, 50, 50, 51, 52, 52, 500, 45, 46, 47, 47, 48, 48, 49, 49



Job to do

Outliers

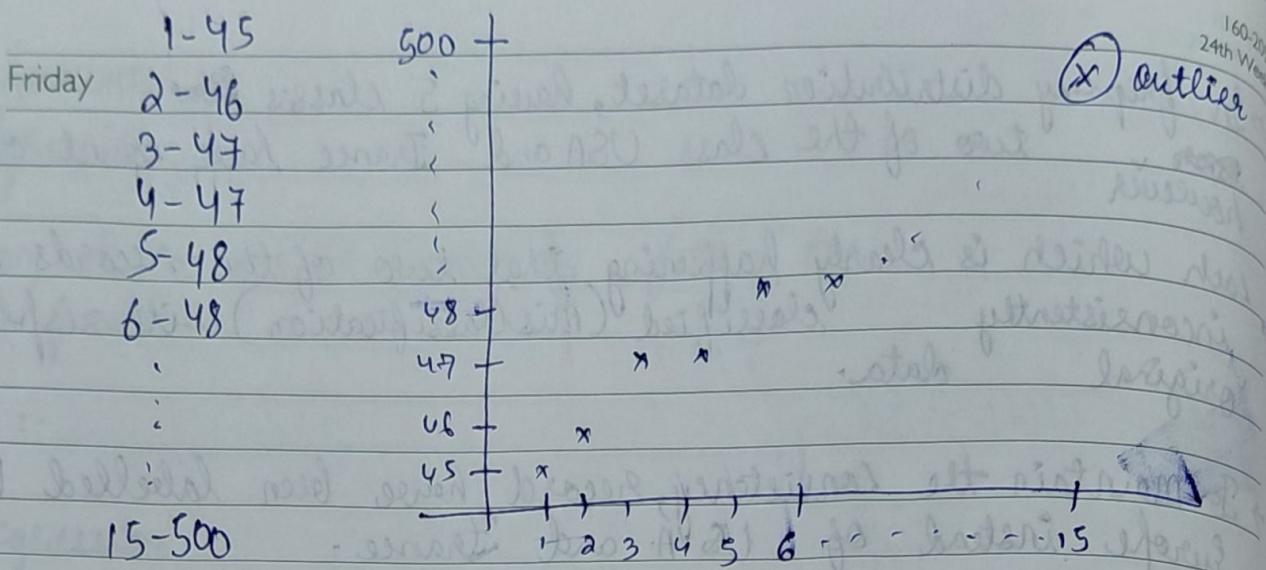
Phone No.



JUNE

2023

09



## Data Integration

→ It is the merging of data from multiple data sources.

It helps to reduce and avoid redundancy and inconsistency in the resulting dataset which improve the accuracy and speed of the ~~subsequent~~ subsequent data mining process.

## Entity Identification Problem

Eg:- Employee Details

Customer Details

Employee Name	Mobile No	Email ID	Customer Name	Mobile No	Email
Rina	99887766 55	Rina@gmail.com	Rina	99887766 55	Rina@gmail.com
Essential	Job to do		Phone No.		

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25

161-204  
24th Week

10

→ It occurs when we try to determine whether two or more records from different data subsets refer to the same real world entity.

Saturday

### Solution

→ Rule based Matching

→ Probabilistic Matching

→ Machine Learning Algorithms

### Data Transformation

→ It is the process of converting data from one format or structure to another format or structure to ensure compatibility, consistency and better quality of analysis.

11 Sunday

### Types

→ Smoothing (Binning: Mean, Median, Boundary), Regression Clustering : It is used to remove the noise from data

→ Aggregation - It summarizes the data. Eg:- Average, Mean, Sum

→ Generalization - Replacing detailed data with higher level concept

Essential

Job to do

Eg:- Age of student  $\rightarrow$  23  $\rightarrow$  Youth Adult

→ Normalization - Rescaling the values to a standard range i.e from 0-1.

12

↳ Attribute Construction: Create a new attribute from the existing data. Eg:- Use DOB to calculate age.

Min - Max Normalization (Scaling) Standard Range - (0-1)

→ It is a linear transformation on the original data

$$X_{\text{nm}} / X^* = \frac{X - \text{Min value of } X}{\text{Range (Max - Min) of } X}$$

where  $X$  is the original data,  $\text{Min of } X$  is minimum value of the attribute,  $\text{Max of } X$  is maximum value of attribute,  $X_{\text{nm}} / X^*$  is normalised value

Q) You have a dataset 50, 60, 70, 80, 90, 100. Min Max Normalisation for 80.

$$X^* = \frac{80 - 50}{100 - 50} = 0.6$$

Q) Minimum and Maximum Value for the attribute Income are 12000 and 98000 respectively. We like to map income to the range 0-1 by min max Normalisation of a value of 73600. Calculate transformed value

$$X^* = \frac{73600 - 12000}{98000 - 12000} = \frac{61600}{86000}$$

$$98000 - 12000$$

$$= 0.716$$

Mid Range = Max Value + Min Value

~~$$\frac{61600}{98000} = 0.628$$~~

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25

164-201  
25th Week

Normalised value of Mid Range

$$X^* = \frac{55000 - 12000}{98000 - 12000} = \frac{43}{86000} = 0.4999999999999999$$

Tuesday

13

Z score Standardisation

$$Z\text{-Score} = \frac{X - \text{Mean}(x)}{\text{SD}(x)}$$

decimal Scaling

Standard Range (-1 to 1)

→ It is a Normalised technique where Normalized value lies between -1 to +1.

$$X_{\text{decimal}} = \frac{X}{10^d}$$

where d = No. of digits in the data value with the largest absolute value

Q) -987, -120, 56, 300, 850 : dataset

$$X_{\text{decimal}} = \frac{-987}{10^3} = -0.987$$

→ -0.987, -0.120, 0.056, 0.300, 0.850

Essential

Job to do

Phone No.

14

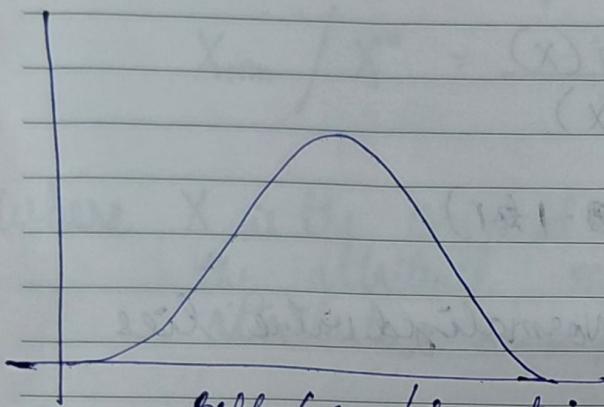
Wednesday

Transformation to achieve Normality

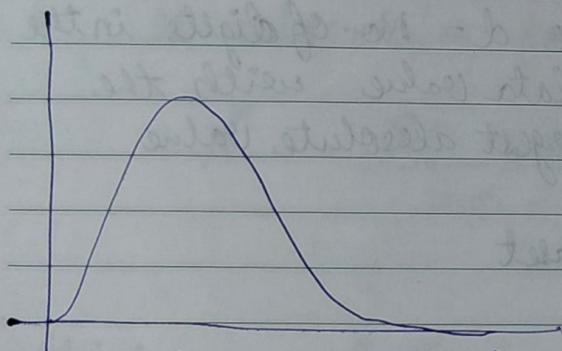
→ Normal distribution is a continuous probability distribution commonly known as Bell curve which is a symmetric curve.

Condition of Bell curve -  $SD = 1$ , Mean = 0

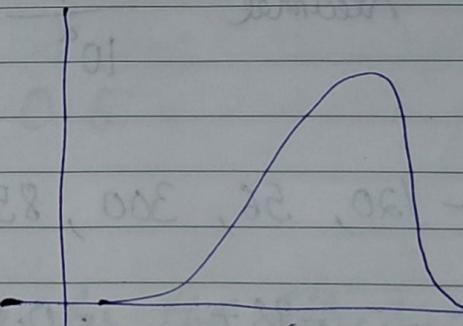
Symmetric curve - Mean, Median, Mode are peak of curve



Bell Curve / Symmetric Curve



+ve (Mean > Median)



-ve (mean < median)

$$\text{Skewness} = \frac{3(\text{Mean} - \text{Median})}{\text{SD}}$$

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17

166-199  
25th Week

15

## Methods to remove skewness

Thursday

1. Log Transformation  $\rightarrow \bar{x} = \log(x+1)$

2. Square Root Transformation  $\rightarrow \bar{x} = \sqrt{x}$

3. Inverse Square Root Transformation  $\rightarrow \bar{x} = \frac{1}{\sqrt{x}}$

## Flag Variables

→ These are dummy / indicator variables which converts categorical variables into only two values (0 or 1).

→ Eg:-

Categorical Variable: Gender (Female, Male)

If gender =: Female then Gender-Flag = 0

If ~~gender~~ gender =: Male then Gender-Flag = 1

Flag =: If region = North then North-flag = 1 otherwise  
North-flag = 0

South-flag =: If region = South then South-flag = 1 otherwise  
South-flag = 0

Essential	Job to do	Phone No.
East-flag =: If region = East then East-flag = 1 otherwise East-flag = 0		

Otherwise All East-flag = 0

16

Friday

Categorical values must be converted into numerical format for machine learning algorithms :-

### ① Label Encoding

Eg:- size

- Large - 0
- Medium - 1
- Small - 2

### ② One Hot Encoding

Eg:- colors

- Red
- Green
- Blue

	Red	Green	Blue
Red	1	0	0
Green	0	1	0
Blue	0	0	1

### ③ Binary Encoding

Eg:- Binary Encoding

- Bhubaneswar 1 → 0001
- Cuttack 2 → 0010
- Rourkela 3 → 0011

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	
16-197 25th Week																	

## DMPA Questions

17

Saturday

- a) Suppose that a hospital tested the age and body fat data for some randomly selected adults.

Age	23	23	27	27	39	41	47	49	50
% fat	9.5	26.5	7.8	17.8	31.4	25.9	27.4	27.2	31.2

a) Calculate Mean, Median, Mode and SD of age and % fat

b) Transform above numeric attribute using Z-score Normalization.

$$\begin{aligned} \text{age} \\ \text{Mean} &= 36.22 \\ \text{Mode} &= 23, 27 \text{ (Bimodal)} \\ \text{Median} &= 39 \end{aligned}$$

$$\begin{aligned} \text{SD} &= \left( \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \right)^{\frac{1}{2}} \\ &= \sqrt{\frac{1}{9} \times 1019.5556} = 10.64 \end{aligned}$$

$$\begin{aligned} \text{fat} \\ \text{Mean} &= 22.74 \\ \text{Mode} &= \text{NonModal} \\ \text{Median} &= 26.5 \end{aligned}$$

$$\begin{aligned} \text{SD} &= \left( \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \right)^{\frac{1}{2}} \\ &= \sqrt{\frac{1}{9} \times 635.2024} \\ &= \sqrt{70.5780} = 8.401 \end{aligned}$$

18 Sunday

Essential

Job to do

Phone No.

JUNE

2023

19

Monday

$$z\text{-score} = \frac{x - \mu}{\sigma}$$

z-score

age	-1.242	-1.242	-0.866	-0.866	0.261	0.449	1.010
% fat	-1.576	0.447	-1.778	-0.588	1.030	0.376	0.554
	1.201	1.295					
	0.530	1.007					

170-19  
26th Week

⑥ Calculate correlation coefficient  $r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{n\sum x^2 - (\sum x)^2}}$

$$= 0.71$$

## A) BINNING

Sales Records :  $T = \{50, 55, 65, 72, 11, 13, 15, 35, 92, 108, 150, 5, 8, 10, 18, 7, 204, 210, 215\}$

Partition them into 4 bins and apply smoothing by bin boundaries

a) Equal width      b) Equal depth

Essential

Job to do

Phone No.

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25

171-194  
26th Week

Sorting: 5, 8, 10, 11, 13, 15, 35, 50, 55, 65, 72, 92, Tuesday  
108, 150, 187, 204, 210, 215

20

min = 5, Max = 215, Range = ~~200~~ 210

$$\text{Width} = \frac{210}{4} = 52.5$$

Bin 1 [5, 57.5] = [5, 8, 10, 11, 13, 15, 35, 50, 55, ~~65, 72, 92~~]

Bin 2 [57.5, 105] = [62, 72, 92]

Bin 3 [105, 155] = [108, 150]

Bin 4 [155, 205] = [187, 204, 210, 215]

Equal Depth ] Take 4, 4, 5, 5

Transforming Categorical Variable into Numerical Variables

e.g.: Region (East, West, North, South)

Categorical Variable	Numeric Value
----------------------	---------------

East

1

West

2

North

3

Essential

South

4

Phone No.

Relation

1 < 2 < 3 < 4

East < West < North < South

21

Wednesday

→ This simply transforms categorical variable region into a single ~~attribute~~<sup>172-193  
26th Week</sup> numerical variable rather than using several flag variables.

### Reclassifying Categorical Variable

- Reclassifying categorical variables means combining regrouping or transforming categorical variables into fewer and more meaningful groups to improve data analysis or model performance.
  - When a categorical variable has too many categories or irrelevant categories, we merge or rename them to make analysis easier.
- Reason -
- Too many categories exist. Some categories have very few observations.
  - Several categories have similar meaning.
  - To create more useful or predictive grouping for data mining algorithms.

Eg:- 50 states can be reclassified on the basis of economic-level as Richest, Middle Range and Poor

Essential	Job to do	Phone No.
	P>E>F>I A>B>C>D>E>F	

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S						
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30

173-192  
26th Week → Adding an Index Field

Thursday

22

- An Index field is a new variable / column added to the dataset to uniquely identify each record.
- It is especially useful when the dataset does not contain unique key or id.  
Why required?
- To uniquely identify each record, data cracking, debugging, joining dataset

eg:-	<u>Index</u>	<u>Name</u>	<u>Age</u>	<u>Salary</u>
	1	Rima	32	20K
	2	Aditya	55	50K
	3	Raj	15	70K
	4	Nakala	22	55K

### Removing variables that are not useful

- Some variables in the dataset may not contain useful information for data mining task, these variables should be identified and removed to simplify the dataset and improve performance of the model.

#### Reason

i) To reduce the noise (irrelevant data)

Essential  
ii) To improve computation

Phone No.

iii) Simplify interpretation

iv) Avoid redundancy

23

Friday

Eg:- Unary Variable: It is a variable/column that has only one value for all the records in a dataset.

Eg:- Table

Index	Name	Age	Salary	Country
1	John	30	50000	India
2	David	32	55000	India
3	Mike	33	60000	India
4	Alex	34	65000	India

Unary Variable

Country
India
India
India
India

Nearly Unary Variable

Index	Name	Age	Salary
1	John	30	50000
2	David	32	55000
3	Mike	33	60000
4	Alex	34	65000

Nearly  
Unary  
Variable

Country
India
India
India

Variables that should probably not be removed

Some variables that appear unimportant at first but later turn out to be quite useful

Types of variables

→ Categorical identifiers that group records meaningful

Essential	Job to do	Phone No.
→ Variables with <del>no</del> potential interaction		effect
→ Variables with missing / <del>rare</del> rare values		

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	M	T	W	T	F	S								
1 175-190 26th Week	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	.

Chapter - 3

## Exploratory Data Analysis (EDA)

Saturday

24

→ There are two distinct approaches to data analysis :-

① Hypothesis Testing :

is confirmatory formal procedure that ~~tests~~ a pre-specified idea or assumption.

② EDA (Exploratory Data Analysis) - It is an open ended discovery oriented process where the role is to learn what the data suggest without a fixed hypothesis.

Both approaches play a complementary role in data mining, statistics and machine learning.

25 Sunday

### HYPOTHESIS TESTING

- It is a statistical method used to make decision / draw conclusion about a population based on sample data.
- It helps to determine whether there is enough evidence to support or reject a particular belief (Hypothesis).

Essential

Job to do

Phone No.

JUNE

2023

26

Monday

Key features :-

i) It starts with an ~~an~~ prior hypothesis or assumption  
(before examining the data in detail.)

ii) It involves two ~~complimenting~~ competing statements

$H_0$  = Null Hypothesis

$H_1$  = Alternative Hypothesis  
or  $H_A$

Example:- A mobile phone operator takes a hypothesis:

$H_0$  = Market share has not decreased after fee hike.

$H_1$  = Market share has decreased after fee hike.

Exploratory Data Analysis

→ It is an approach to analyse dataset that emphasises ~~weak~~ exploration and descriptive statistics to uncover patterns ~~&~~, anomalies and relationships without predetermining assumptions.

visual

Essential Objectives

Job to do

Phone No.

i) Understanding the structures of data - e.g:- Variables, Range

ii) ~~Histograms~~ distribution of categorical variables

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25

178-187  
27th Week

27

- iii) Look at the histograms of numeric variables to understand their spread and shape.
- iv) Explore the relationship among sets of variables (Predictor ~~set~~ and Target Variables)
- v) Detect Outliers, missing values and data quality issues.
- vi) Develop initial hypothesis and guide subsequent modelling.

### Common EDA techniques

- Graphical - Histogram, Scatter Plot, Box Plot, - . -
- Numerical - Summary statistics (mean, median, mode, variance, skewness), correlation coefficient
- Subset / group analysis -

Identifying clusters, ~~sets and~~ trends or interesting subsets.

- EDA acts as foundation of data analysis shaping the direction of further investigation and hypothesis testing.

Essential

Job to do

Phone No.

28

## Wednesday Complementary Roles

→ EDA comes first (discovery stage)

↳ Helps analyst understand the dataset, distribution and uncover important relationships and patterns that could indicate important areas for further investigation.

→ Hypothesis testing follows (confirmation stage)

↳ Validates the patterns or suspicions suggested by EDA with statistical reasons i.e. testing assumptions with formal procedures.

→ Together they form a powerful cycle of discovery and confirmation in Data Mining and statistical analysis.

	Hypothesis Testing	EDA
Purpose →	To perform or reject a pre-specified	→ To discover patterns, understand distributions and generate new ideas.
Approach →	Deductive, confirmatory	→ Inductive, discovery oriented
When used →	When clearer theory driven questions exist	→ When data are unfamiliar, large or complex
Focuses →	Formal Decision Making	→ Investigation of variables, distributions and relationships

T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25

180-185  
27th Week

29

Tools → Statistical tests  
(ANOVA, Chi-square)

Graphical (Histogram, Thursday  
Scatter plot, Box Plot)

Outcome → Binary decision

→ Insights, hypothesis

Flexibility → Rigid, structured

→ Flexible, iterative

### EDA on the Churn dataset (A case study)

→ The churn dataset (UCI ML repository) is used to demonstrate EDA methods applied in a real world business scenario.

→ EDA helps in:

- ↳ Detecting anomalies or missing data
- ↳ Identifying patterns and relationships among variables
- ↳ Suggesting potential predictions for target variable
- ↳ Gaining domain insights through visualisation

### Overview of dataset

Essential → No. of observations <sup>to do</sup> = 3333

Phone No.

No. of Inputs/predictors/features = 20

Target (Churn) - Indicates whether a customer has left the company.

JUN / JUL

2023

30

→ The dataset contains a mix of categorical, integer value and continuous features describing customer demographics, account information, service usage, etc.

### Variables in Churn Dataset

#### a) Customer Identification

→ State : Categorical

→ Account Length : Integer

→ Area Code : Categorical

→ Phone Number : Unique Identifier

#### b) Service Plans

→ International Plans : Dichotomous Categorical (yes/no)

→ Voicemail Plans : Dichotomous Categorical (yes/no)

→ No. of voice mail messages : Integer

#### c) Usage Metrics

→ Total minutes, calls, charges, ... (continuous Integer)

#### d) Customer Service Interaction

→ No. of calls to customer service (Integer)

S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

182-183  
27th Week

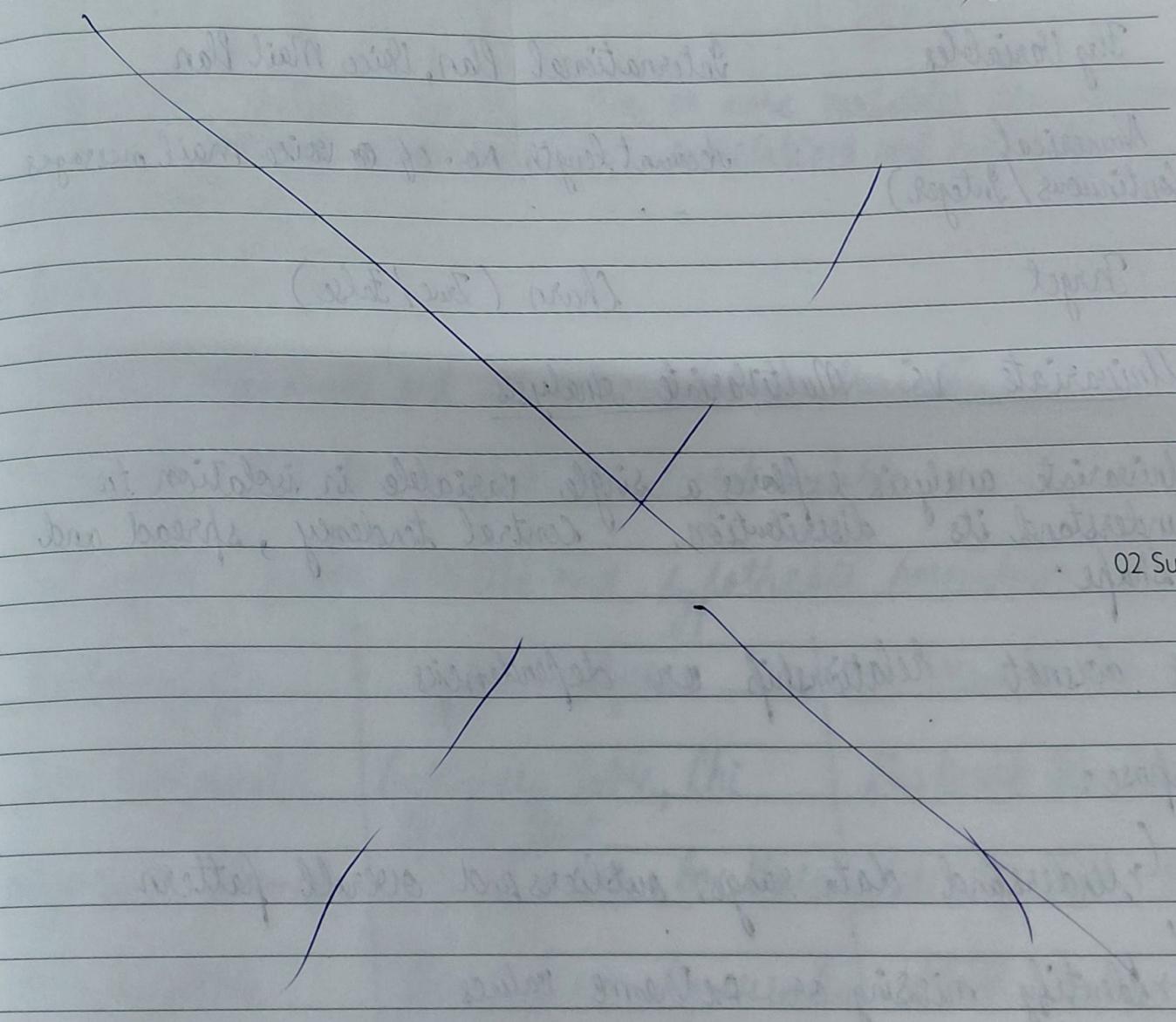
e) Target Variable,

Saturday

01

→ Churn - Boolean (True / False)

indicates whether a customer left the company



02 Su

Essential	Job to do	Phone No.

JULY

2023  
184  
28th W

03

Type Monday

Variables

Categorical

State, Area Code

Identification

Phone Number

Flag Variables

International Plan, Voice Mail Plan

Numerical  
(continuous / Integer)

Account length, No. of voice mail messages

Target

Churn (True / False)

### Univariate vs Multivariate Analysis

- Univariate analysis explores a single variable in isolation to understand its distribution, central tendency, spread and shape.
- It does not relationship or dependencies
- Purpose:
  - ↳ Understand data range, outliers and overall pattern
  - ↳ Identify missing or extreme values
  - ↳ Decide on data transformations ('Normalisation and Log Based Transformation')
  - ↳ Check assumptions

S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M														
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

185-180  
28th Week

# 04

## Type of Variable

## Common Techniques

## Visualisation

Categorical

Frequency counts, mode, ~~Box Plot~~, Bar chart, ...

Numerical

Mean, median, ... Box Plot, Scatter plot

→ Multivariate analysis investigates two or more variables simultaneously to detect patterns, relationships, correlations and interactions between them.

→ Purpose:

- ↳ Find dependencies and interactions between variables.
- ↳ Identify predictors (inputs) for target variable.
- ↳ Support future selection and hypothesis formulation

## Relationship Type

## Typical analysis

## Visualisation

Two Categorical

Contingency table, Chi square test

Clustered Bar Chart

1 categorical + 1 numeric

group means, box plots

Side by Side Box Plot

Two numeric

Correlation, regression line

Scatter Plots

Essential Many numeric

Job to do Heatmap

Phone No. Matrix Plots

JULY

202

05

Wednesday

## Contingency Table

→ It is also known as Cross Table

→ It is a type of table using statistics to show the frequency distribution of the variables i.e. how two or more categorical variables are related to each other.

Job to do

Phone No.

