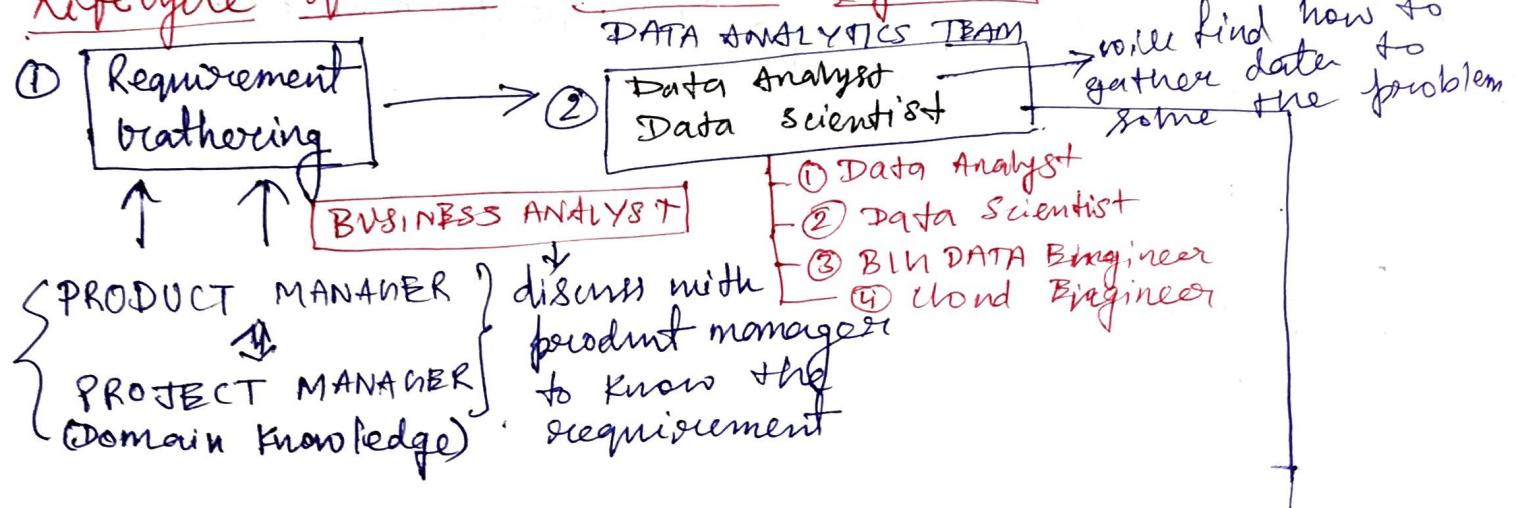
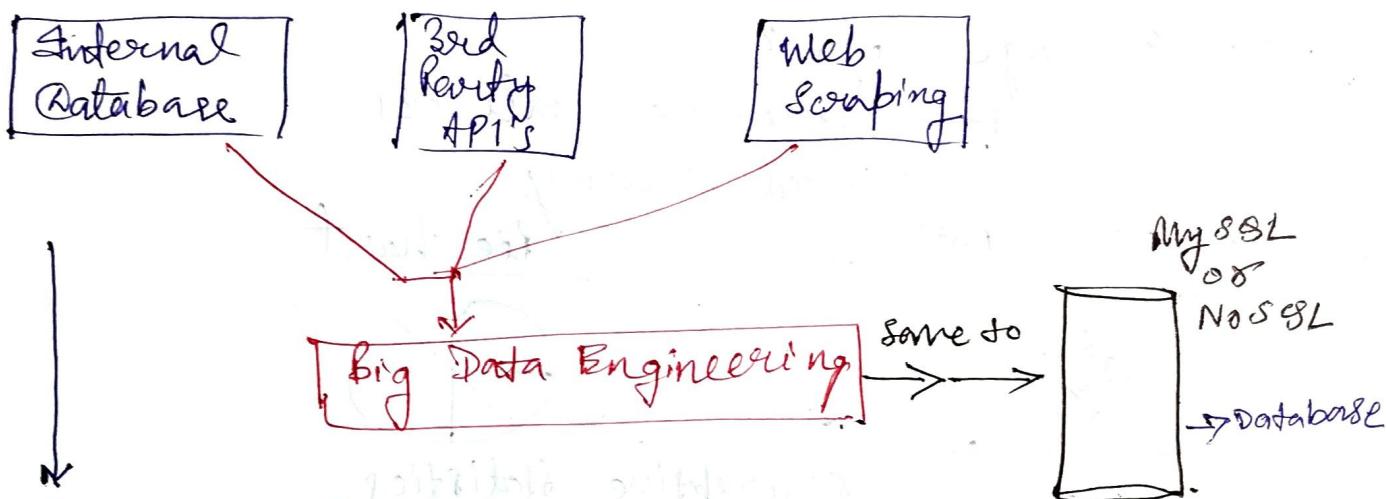


Lifecycle of a DATA Science Projects.:

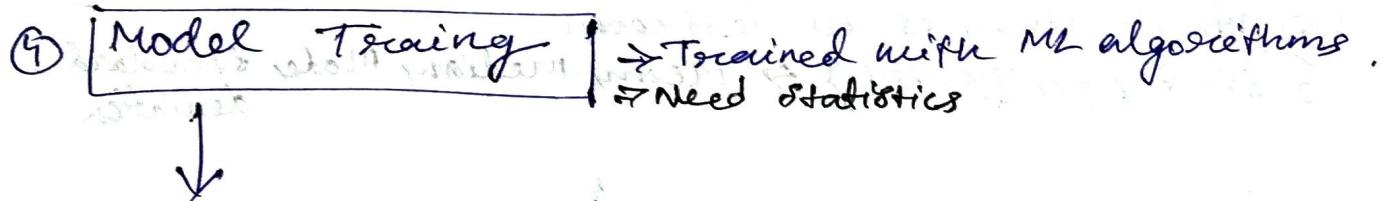
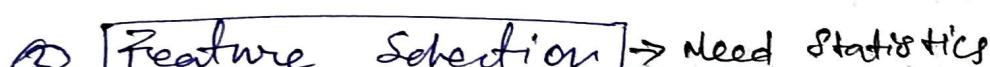
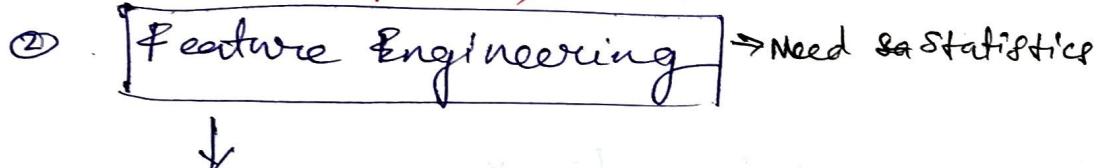
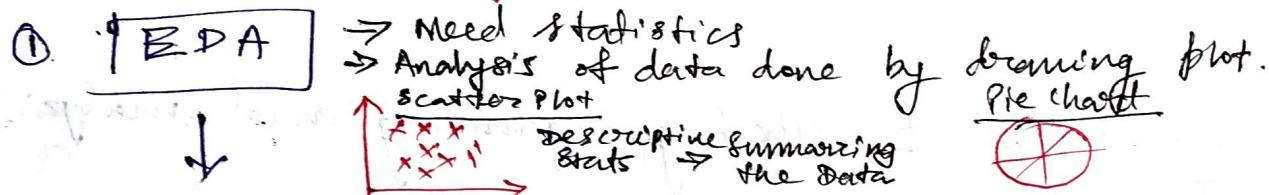


③ Sources of Data



④ Data Science Projects

→ Data Scientist works starts here



⑤ Hypertuning → Improve the performance of the model

⑥ Deployment

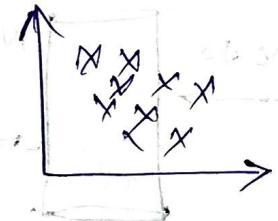
Analysis of Data

① Age = {12, 13, 14, 18, 20, 25}

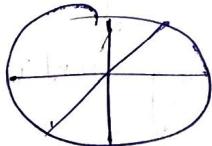
Find average age?

→ This is also descriptive statistics.
Measure of central tendency.

② Scatter Plot



Pie chart



Descriptive Statistics

Statistics

→ It is the science of collecting, organizing and analyzing of data.

Data

→ "Parts are pieces of information".

e.g. ① Ages of students in classroom
 $\{24, 25, 28, 29, 28\} \Rightarrow$ Mean, Median, Mode, Standard deviation

② Weights of students in classroom

$\{85, 65, 72, 73, 89\} \Rightarrow$ Mean, Median, Mode, Standard deviation

Types of Statistics

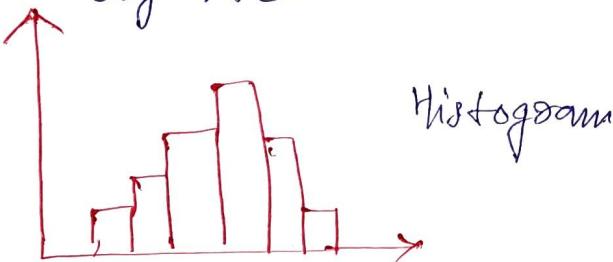
- ① Descriptive Statistics
- ② Inferential Statistics

Descriptive Statistics

→ It consists of organizing and summarizing of data using different plots.

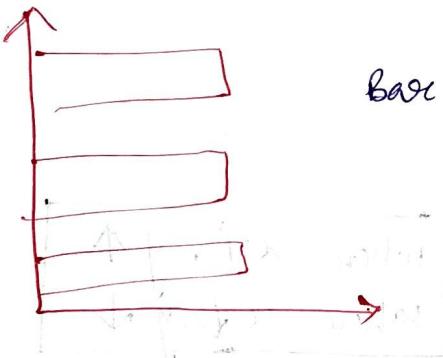
→ Extensively used in BPA and Feature engineering.

Eg: ①



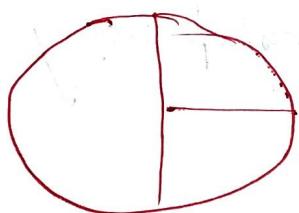
Histogram

②



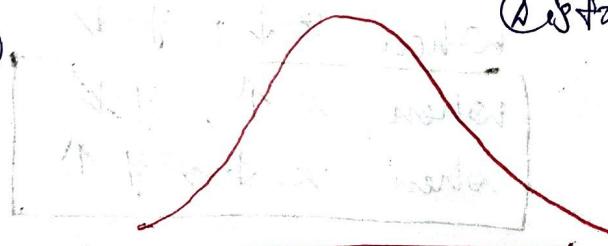
Bar chart

③



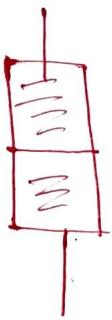
Pie chart

④



Distribution

⑤



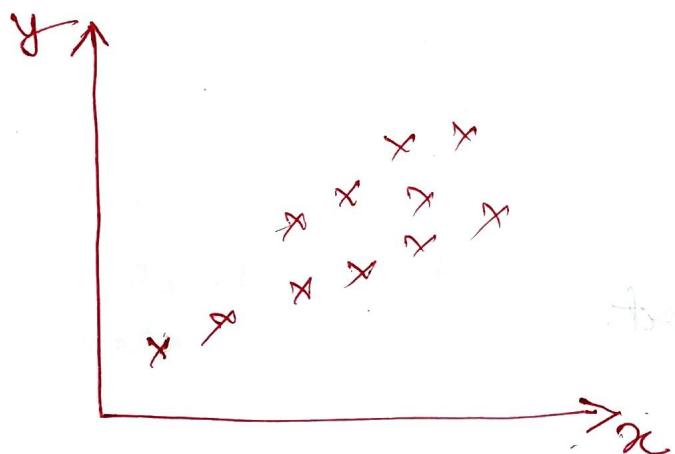
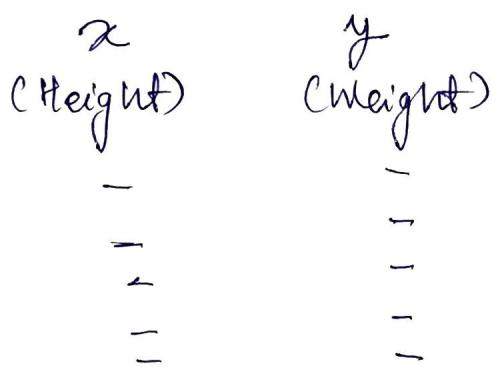
Candlestick

⑥



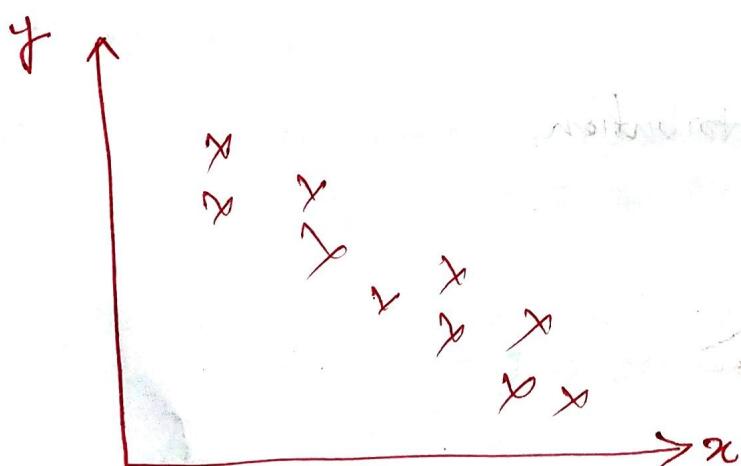
Box Plot

⑦



when $x \uparrow, y \uparrow$
when $x \downarrow, y \downarrow$

→ We can check this relationship using scatter plot.



when $x \downarrow, y \uparrow$
when $x \uparrow, y \downarrow$
when $x \downarrow, y \uparrow$

Inferential Statistics

- It consists of collecting sample data and making conclusion about population data using some experiments.
- Making conclusion → can be done by hypothesis testing.

Eg: ① University → 500 students

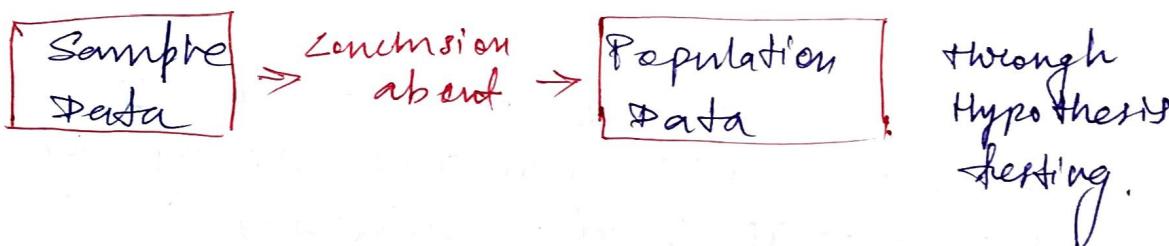
Class A → 60 students



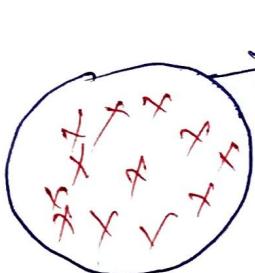
Sample data → Age → Average age of the entire university.

Can we make con average age of the university with the help of sample data of class A?

→ Yes, we can make the prediction.



Sample Data vs Population Data



Punjab
Total Population ≥ 10 cr

Exit Poll.

→ They have to take sample of data.

Sample Size = 1000

Take sample size from every region and then do average and comes up with exit poll.

- PARTY A will win (Earlier/Prediction)
- PARTY B will loose. (Actual result)
- Prediction goes wrong, hypothesis testing goes wrong.

Problem:

Let's say there are 20 classrooms in a university and you have collected the age of students in one classroom?

Ages = {21, 20, 18, 34, 17, 22, 24, 25, 26, 23, 22}

Weights = { - - - - - }

Descriptive Stats:

- Ques: ① What is the average age of students in the classroom?
 ② Relationship b/w age and weight?

Inferential stats:

Ques: ① Are the average of the students in the classroom less than the average age of the students in the university?
 Comparing Sample and Population data.

②

1000 → Students

Class A → 50 girls → 95% marks

Class B → 50 boys → 92.5% marks

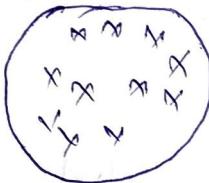
Are the average marks of by girls in the university greater than by boys?

→ We can do this using Hypothesis testing.

Sampling Techniques

① Simple Random Sampling

→ Every member of the population (N) has an equal chance of being selected for your sample (n).



→ We can perform random sampling for Exit Poll, General Survey etc.

② Stratified Sampling

→ Strata → Means layers → Clusters, → barangay
→ barangay the population and then select.

bender → Make

bender → Female

Education Degree →

- High School
- Master
- PhD

Blood groups →

- A+
- B+
- O+

E.g.: Exit Poll \rightarrow age ≥ 18 and age < 18 (Count vote)

③ Systematic Sampling

AIRPORT }

CREDIT CARD sale place

Every 8th person approach \leftrightarrow Every 9th person approach

→ Select every 9th individual out of population.

④ Convenience Sampling

→ only those who are interested in the survey, will only participate.

e.g.: ① "x" people is interested in Data Science Survey.
We will send survey documents to x people only.

② New job forms → fill the form (Interested candidates)

Ques.

① Survey regarding new technology?

→ Convenience Sampling.

② RBP Survey for women to know the expenses?

→ Stratified Sampling, then Random Sampling for married women.

③ Credit card calls?

→ Stratified Sampling and then random sampling.
(for balanced portion)

Variable

→ It is a property that can take any value.

e.g.: age = 14

age = 33

age = 100

Variables

ages = [24, 25, 26, 27, 28, 29] ⇒ collection

Types of Variable:

- ① Quantitative Variables
- ② Qualitative Variables
 - ↳ Categorical Variables

Discrete variable

Continuous variable

Quantitative Variable

- Measured numerically & mathematical operations
- e.g.: Age, weight, height, rainfall (cm), temp., distance.

Qualitative Variable

- Categorical Variables & Based on some characteristics
↳ they are grouped together.
- e.g.: Gender, Types of flowers, Type of movies

Quantitative Variable

① Discrete Variable

- whole number.

- e.g.: ① Number of Bank account
- ② number of children in a family

② Continuous Variable

- continuous, decimal

- e.g.: ① Height
- ② Weight
- ③ Rainfall
- ④ Speed

Assessment

- ① what kind of variable is Marital Status?
→ **Categorical variable**
- ② what kind of variable is Range River Length?
→ **continuous**
- ③ what kind of variable is Movie duration?
→ **continuous**
- ④ What kind of variable is Person?
→ **Discrete** (Person has many varieties).
That's why it is not categorical.
- ⑤ what kind of variable is Pg?
→ **continuous** (Pg doesn't have i-me' names)
- ⑥ what kind of variable is Gender?
→ **Categorical** (Gender has limited categories)

[→ Note: Many number of categories → put it into discrete]

- ⑦ what kind of variable is Pancaid number?
→ **Categorical** (Pancaid number has set of characters)

Pancaid has alphanumeric characters → **Categorical feature**
and it is also a unique identifier