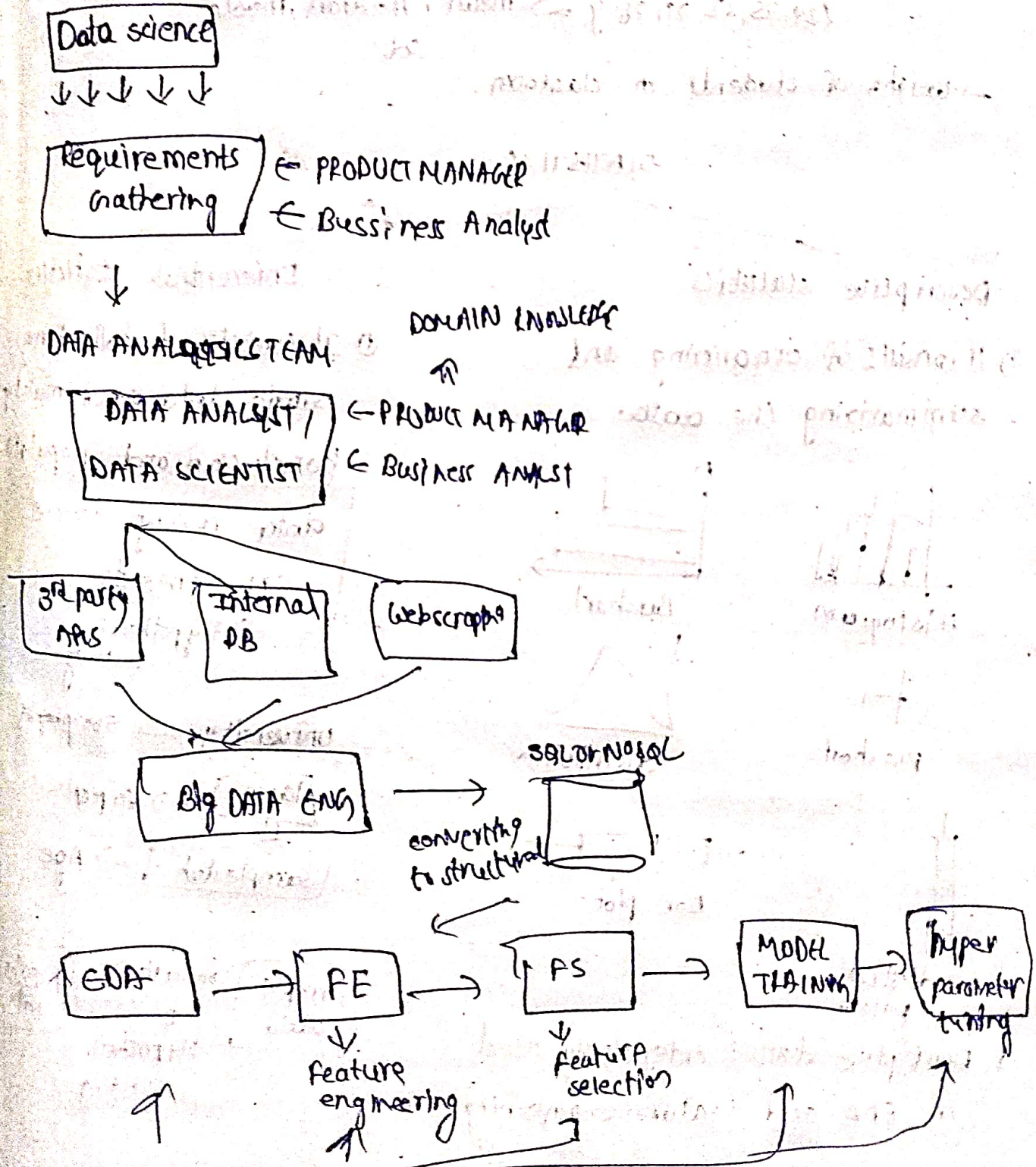


STATISTICS

- * Find the avg size of shark throughout the world?
- * Amazon Big Billion day sale best date? → which month should you select?

LIFE CYCLE OF DATA SCIENCE PROJECT



statistics will be used → Analysis of data and summarizing of data.

STATISTICS →

Statistics is the science of collecting, organizing and analyzing the data.

Data

→ Facts or pieces of information.

Eg: Ages of students in classroom.

{24, 25, 32, 21, 28} → mean, median, mode

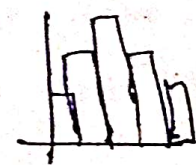
→ weights of students in classroom

sd

STATISTICS

Descriptive statistics

① It consists of organizing and summarizing the data.



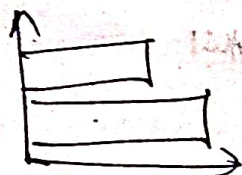
Histogram



Pie chart



Candlestick pattern



Bar chart



Distribution



Box plot

Inferential statistics

① It consists of collecting sample data and making

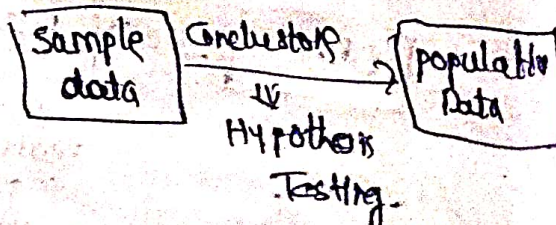
conclusion and population data using some experiments.

→ Hypothesis testing.

University → 500 people

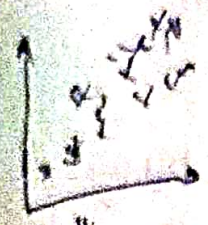
Class → 60 people

Sample data ⇒ Age.



* Descriptive stats is extensively used in EDA and feature engineering

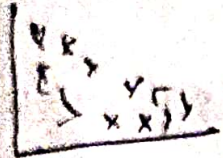
SAMPLE DATA vs population data



scatter plot

x ↑ y ↓

x ↓ y ↑



Sample size = 1000

Punjab

1000
↓
Population DATA

population data (N)

sample data (n)

eg: 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 }

weight { - - - - - }

Descriptive stats

- 1) what is the avg age of students in the classroom?
- 2) Relationship b/w age & weight?

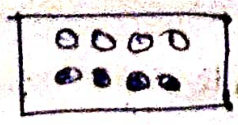
Inferential stats

- 1) Are the avg age of students in the classroom less than avg age of students in the university?

Sampling Techniques

1) Simple Random sampling: Every member of population has an equal chance of being selected for your sample (n).

N = 8



n = 4

picking each marble probability is same.

2) stratified sampling

Sample will be selected based on groups, layers and clusters.

Gender → Male Female

education degree → High school Master PhD

3) systematic sampling \rightarrow Airport n^{th} person

{credit card}

every
3rd person

\rightarrow

every n^{th} person

\rightarrow select every n^{th} individual out of population (N)

④ convenience sampling

only those who are interested in the survey will only be participated.

Ex: AT survey \rightarrow General AT survey

Variable

A variable is a property that can take any values.

Ex: age = 14

age = 25

age = 100

\rightarrow two types of variables:

① Quantitative variable \rightarrow measured Numerically

{mathematical operation}

Ex: age, weight, rainfall, temperature

② Qualitative variable \rightarrow categorical variables {based on some characteristics}

Ex: gender, types of movies

Quantitative variable

\swarrow

Discrete variable \rightarrow

Ex: \rightarrow No. of children \rightarrow whole number

\rightarrow limited

\rightarrow No. of bank accounts

\searrow continuous variable

Ex: age, height, speed, rainfall.

Assessment 1: 100, 150, 180, 200, 220, 250

* What kind of variable is Marital status → Categorical variable

" " " " " Chang river length & continuous

Movie Duration & Continues

pincode \leftarrow Discrete

IA: \rightarrow Continuous

Day-2

AGENDA

→ Histograms

Measure of CT

MEASURE OF DISPERSION

PERCENTILES & QUANTILES

5 NUMBER SUMMARY (BOX PLOT)

Histogram &

Agest = 10, 12, 14, 16, 24, 26, 30, 35, 38, 39, 40, 44, 47, 50, 51, 56, 76, 99, 100

steps to create histogram

① sort the numbers

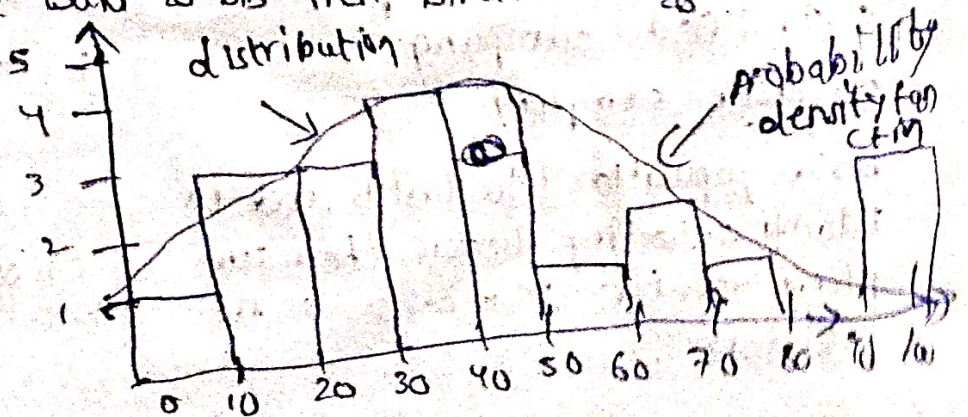
② Bins \rightarrow number of groups

⑧ Bin size \rightarrow size of Bins.

$$\text{Number of Bins} \approx \frac{\text{Max value}}{\text{size of Bin}}$$

size \rightarrow size of Bins.
suppose we want 10 bins then Number of groups = $\frac{100-0}{10} = 10$

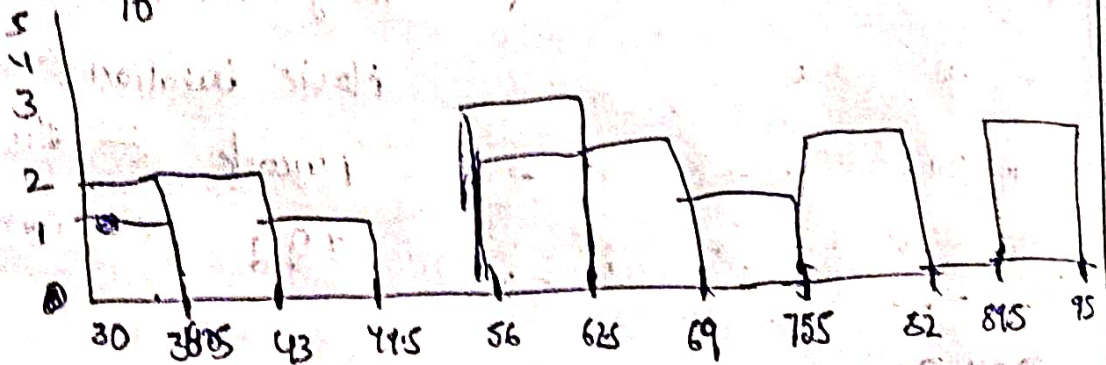
suppose we want 20 bls then blnrite = $\frac{100-0}{20} = 5$



Weight = {30, 35, 40, 42, 46, 50, 54, 62, 63, 68, 75, 77, 80, 89, 95}

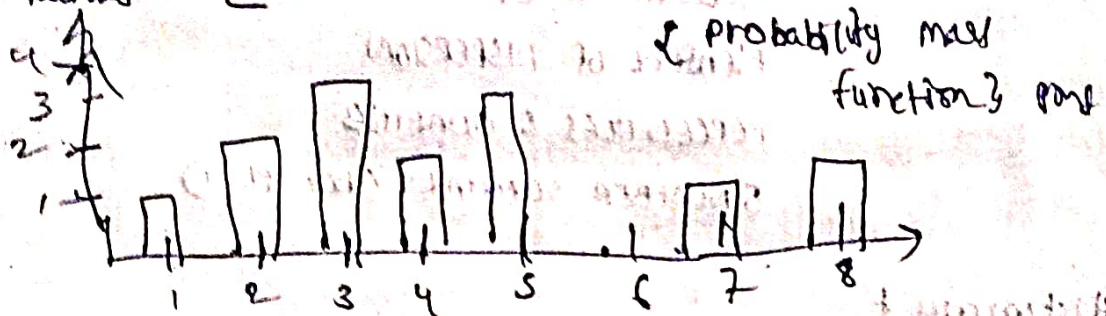
bins = 10

$$\text{bin size} = \frac{95 - 30}{10} = 6.5$$



Discrete Bank account

No. of Banks = [2, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 4, 5]



→ Pdf for continuous variable

→ PMF for discrete variable

Sampling

probability sampling

- simple random sampling
- systematic sampling
- stratified sampling
- cluster sampling

Divide population into clusters and randomly selecting clusters, then sampling all the members in the clusters.

Non probability sampling

- convenience sampling
selecting individual who are easy to reach
- judgemental sampling
select individual based on researcher's judgement
- snowball sampling
existing study subjects recruit future subjects

multi-stage sampling.
combining multiple sampling
method. involves selecting
clusters and randomly
sampling

among their advantages

(b) Quota sampling

Age, Group, Gender, Caste,

→ selecting sampling depending on uses/cases.

Qualitative data

Nominal

↓ No rank

Ex: M/F → Gender

Blood group

Pincode

→ categorical values

Ordinal

= categorical
data w/o

can assign
rank

Good 'Bad better

1 2 3
ext customer
feedback

scales of measurement

→ scales of measurement describe the nature of information within the values assigned to variables

4 primary scales of measurement

1) Nominal scale

2) Ordinal scale

3) Interval

4) Ratio.

this scale classifies data into distinct categories that do not have an intrinsic order.

characteristics

i) Data is categorized based on names, labels, or qualities

ii) categories are mutually exclusive.

iii) No logical order among categories

ordinal scale

this scale classifies data into categories that can be ranked or ordered.

characteristics

- i) Data is categorized and ranked in specific order
- ii) The interval between ranks are not necessarily equal.

eg: customer feedback

satisfied very satisfied not satisfied

Interval scale

The interval scale not only categorizes and orders but also specify the exact difference between intervals. It lacks a true zero point.

characteristics

- 1) Data is ordered with consistent interval between values.
- 2) Allows for meaningful comparison of values.
- 3) No true zero point.

eg: Temp in Fahrenheit
10°F, 20°F, 30°F

calendar years
2020, 2016, 2020

Ratio scale

The order matters. The differences are measurable. Contains a zero starting point.

eg: student marks in a class

ASSIGNMENT#

- ① length of different rivers in world → Interval scale.
- ② favorite food based on gender? → ordinal scale.
- ③ Marital status? → Nominal scale.
- ④ IQ measurement → Interval scale.