

Statistics For Data Science

03/09/2022

1. Use Case :-

Bank → Let suppose HDFC
have already two ATM at
location A & B

(A)

(C)

(B)

This task is basically for
Data Analyst or Data
Scientist.

Q HDFC have already
two ATM at
location A & B

They are planning
whether they will
open new ATM
at location C or
not?

2. Use Case :-

Q. Find the Average size of shark throughout the world?

3. Use Case :-

Q. Amazon Big Billion Day Sale

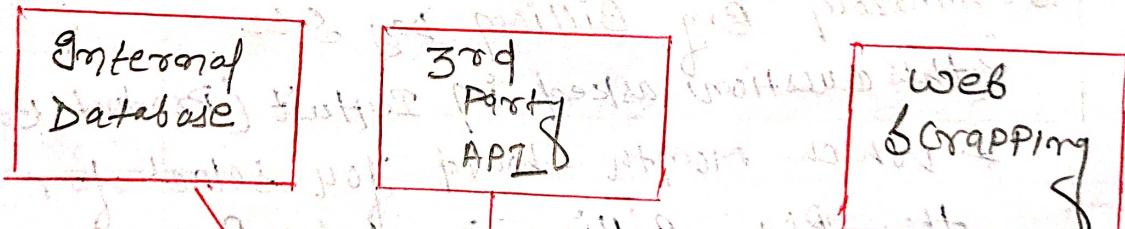
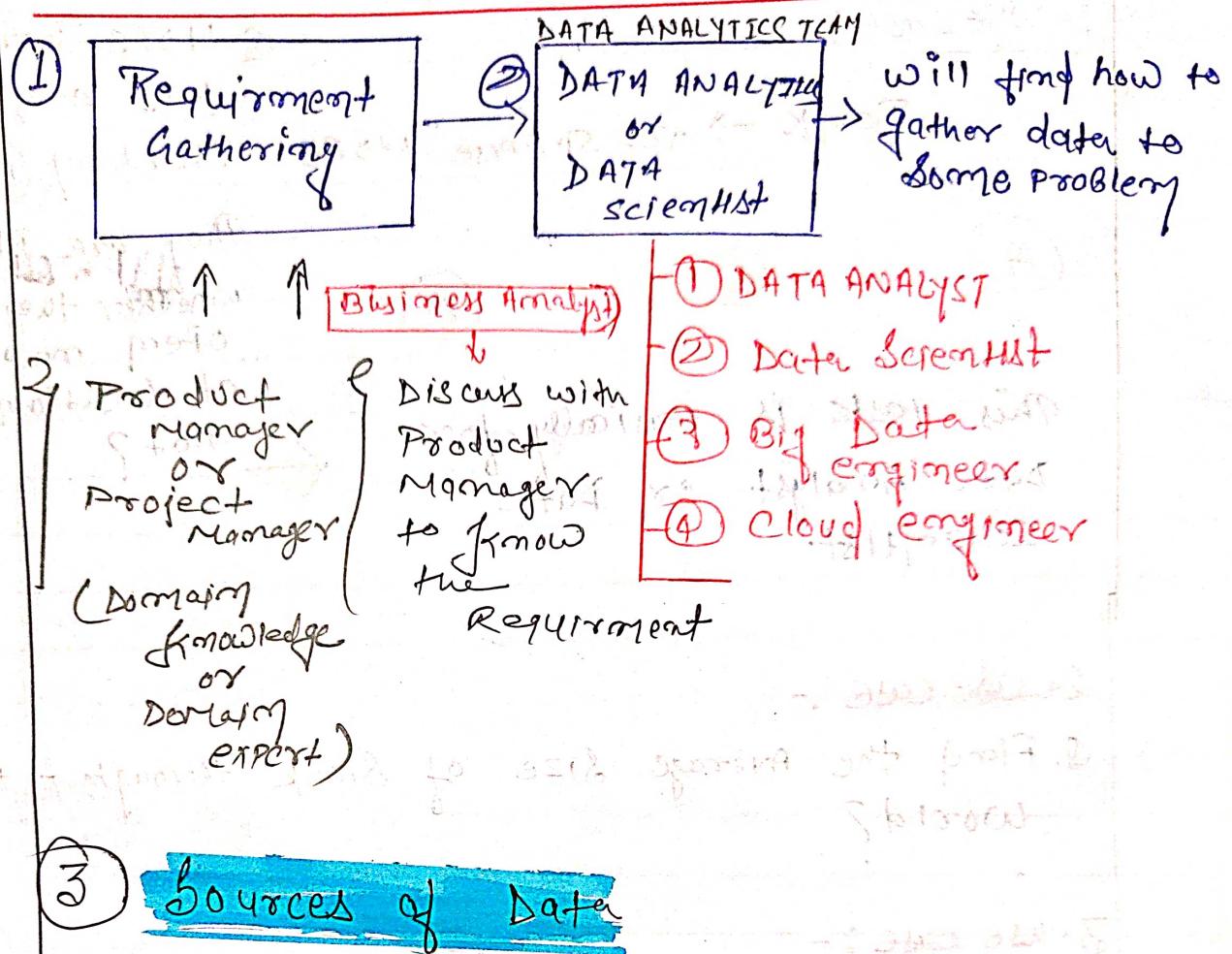
{this question asked in Pehuit (Product Based) company}

* Which Month should you select for
the Big Billion Day Sale?

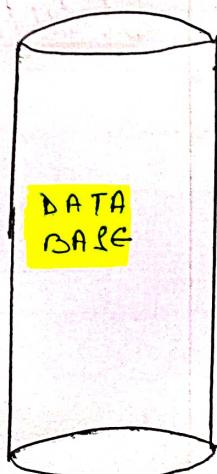
Statistics

[Life cycle of Data Science Project]

Life cycle of a Data Science Project →



Big Data Engineering



My-S91

⇒ we use to save text in my S91.

No S91

⇒ For saving the images we use ~~No S91~~

④

DATA SCIENCE PROJECTS

⇒ The work of DATA Scientist start here.

(i)

EDA :-

Exploratory Data Analysis

(need Statistics)

(ii)

Feature engineering

(iii)

Feature selection

(iv)

Model selection

→ Trained with ML Algo

(v)

Hyper Parameter tuning

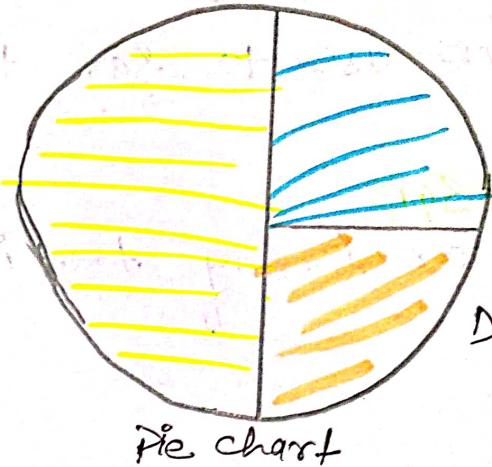
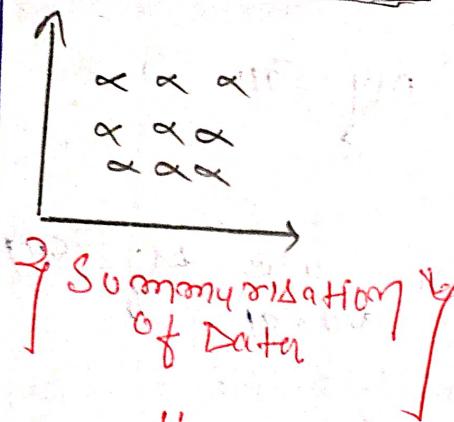
→ Improve the Performance of Model

(vi)

Deployment

* In All the steps Statistics will be used !!

Analysis of Data



Descriptive Stats.

Age = {12, 13, 14, 18, 20, 25}

Average age } is part of Descriptive Stats } is
a measure of central tendency

Interview Question

Q. How Statistics is used in Machine Learning?

Definition

Statistics

→ Statistics is science of collecting & organising & analysing the data.

Data

→ Data is nothing but fact or piece of information

Eg:-

① Ages of students in classroom

{ 24, 25, 32, 29, 18 }

With the help of this data we can analyse this data by finding

Mean, Median, Mode & Standard deviation

Eg 2

weight of students in a classroom

Types of Statistics



Descriptive stats

→ It consists

& summarizing of organization of Data

→ Extensively used in EDA

& feature engineering

inferential stats

⇒ It consists of

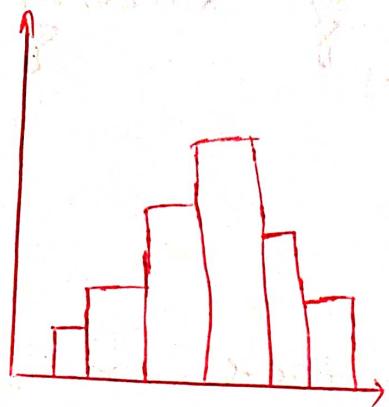
collecting sample

Data & making

Conclusion about Population

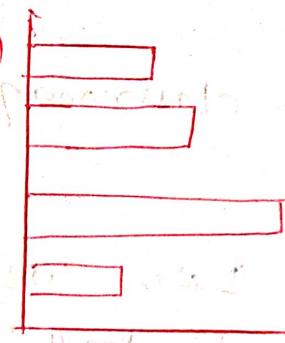
Descriptive Statistics

①



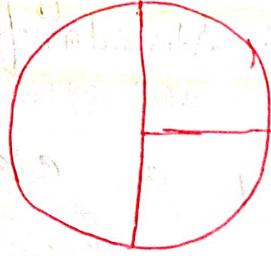
Histogram

②



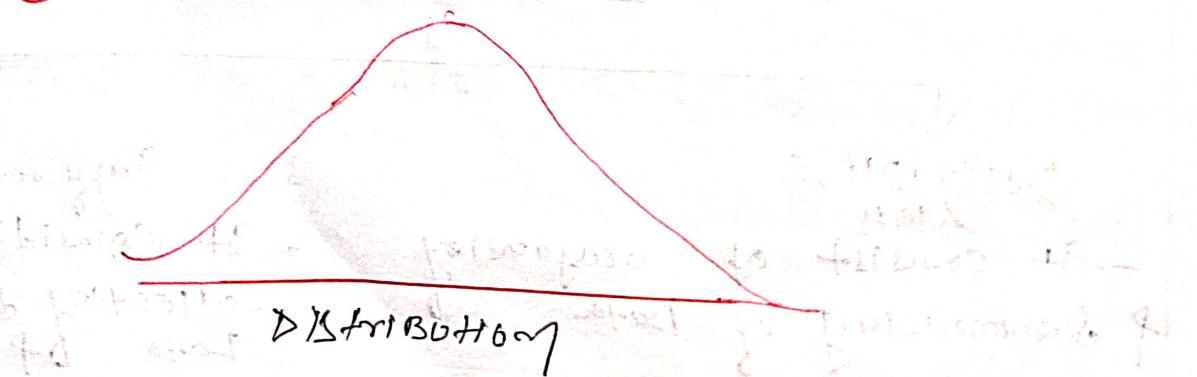
BAR CHART

③



Pie chart

④



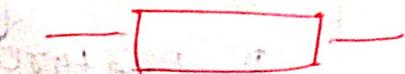
Distribution

⑤



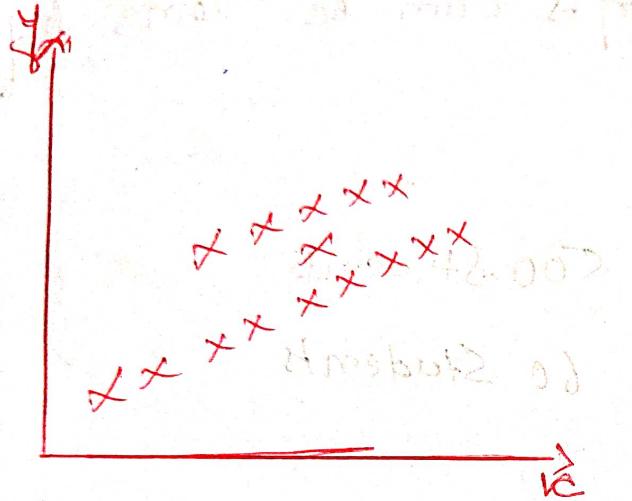
Candle stick

⑥



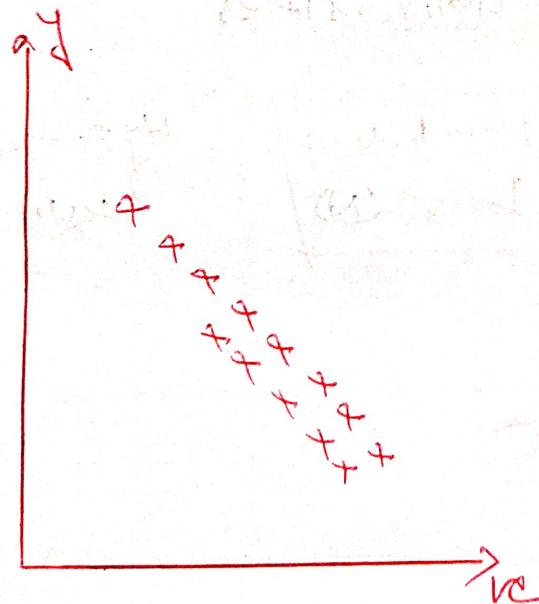
Box plot

| r^2 | Height | Weight |
|-------|--------|--------|
| — | — | — |
| — | — | — |
| — | — | — |
| — | — | — |
| — | — | — |
| — | — | — |
| — | — | — |
| — | — | — |
| — | — | — |



when $r^2 \uparrow$, $y \uparrow$
when $r^2 \downarrow$, $y \downarrow$

→ we can check this
Relationship using
Scatter Plot

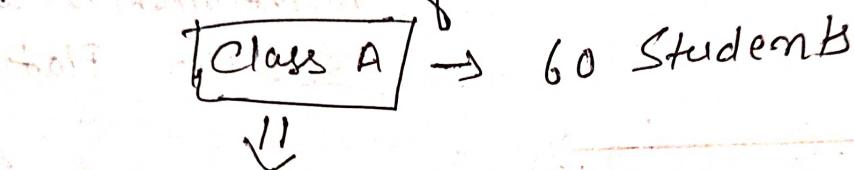


when $r^2 \uparrow$, $y \downarrow$
when $r^2 \downarrow$, $y \uparrow$

Inferential Statistics

- It consist of collecting sample data & making conclusion about population data using some experiment
- Making conclusion → can be done by hypothesis testing

Eg: ① University 500 students



Sample data \Rightarrow Age \Rightarrow Average age of entire university

Sample
Data (m)

CI
Confidence
interval

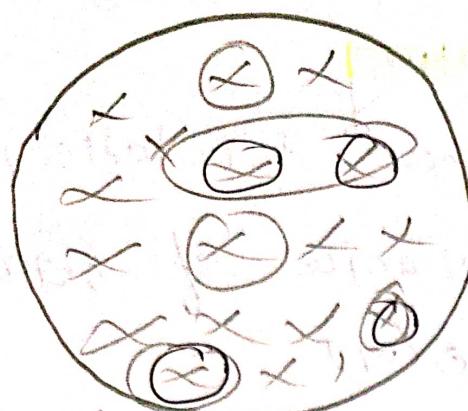
Population
Data (N)

Hypothesis
testing

P-value

- ① Z-test
- ② t-test
- ③ Chi-Square test
- ④ F-test

Sample Data vs Population Data



↑
Pumpaf State
(10 Cr. Population)

↓
Population Data

Needless will conduct some exit poll

We can't go to every people to ask, so we will select some sample of size 1000 & ask from that peoples

Eg: → Let's say there are 20 classrooms in University & you have collected the age of students in one classroom.

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

Weight = ?

?

Descriptive Stats

Q. What is Average age of student in classroom?

Q. Relationship Between Age & weight

Different types of stats

Q. Are the average age of the student in the classroom less than average age of the student in the university?

Population Data = (N)

Sample Data = (n)

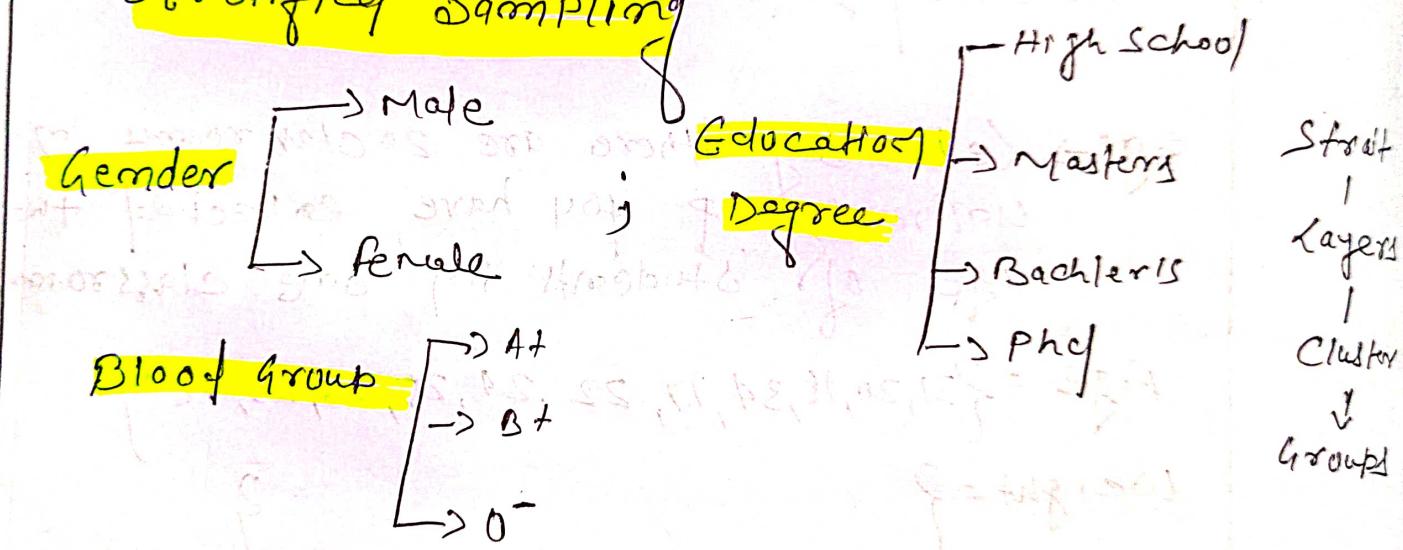
Q. What are the different sampling techniques?

1. Simple Random Sampling

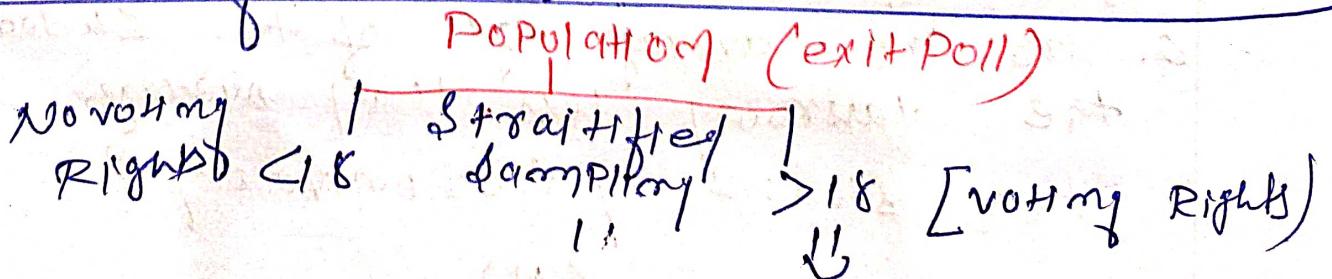
→ Every member of the population (N) has equal no. of chance of being selected for sample (n)

Eg: → exit poll, movie review, lottery

2. Stratified Sampling



In exit poll, we first apply stratified sampling for megalection age ≥ 18 & choosing the age > 18

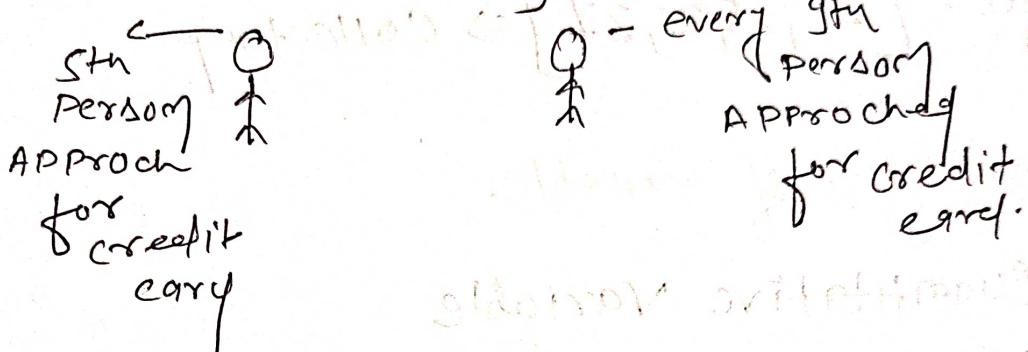


Applying Random Sampling

③ Systematic Sampling

⇒ Selection of every m^{th} individual out of Population (N):

Example: ⇒ At MALL }



④ Convenience Sampling

⇒ Only those who are interested in survey. Only they will participate

Example:-

- ① Survey Regarding my New technology
- ② Credit card calls.
- ③ RBI Survey for know the house expense
Only for women —
Married [→ Stratified
unmarried ↓]
convenience sampling

Variable

→ A variable is property that can take any values
eg: $\rightarrow \text{age} = 14, 25, 100$

Variables

Ages = $[24, 25, 26, 27]$ → collection

Two Types of variable

① Quantitative Variable

→ Measured Numerically

(Mathematical operations of)

eg: \rightarrow Age, weight, Rainfall (cm), temperature, distance

② Qualitative Variables

Eg:- Gender, Type of flower, Type of movie

* Based on some characteristics they are grouped together

Quantitative Variable

Discrete (Whole No) Variable

- * No. of Bank Account
- * No. of children
- * No. decimal value

Continuous Variable

- Eg:- Continuous
- ① Height, weight, Age, Rainfall, Speed.