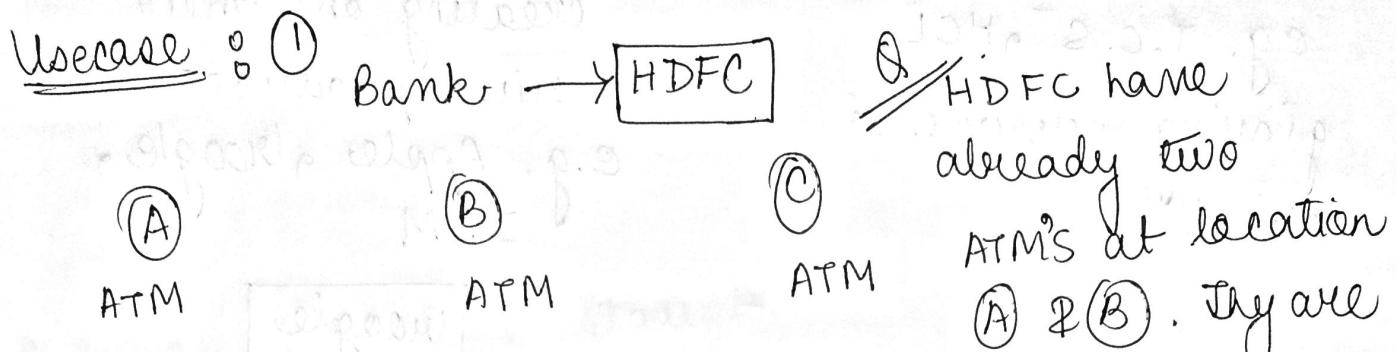


STATISTICS FOR DATA SCIENCE



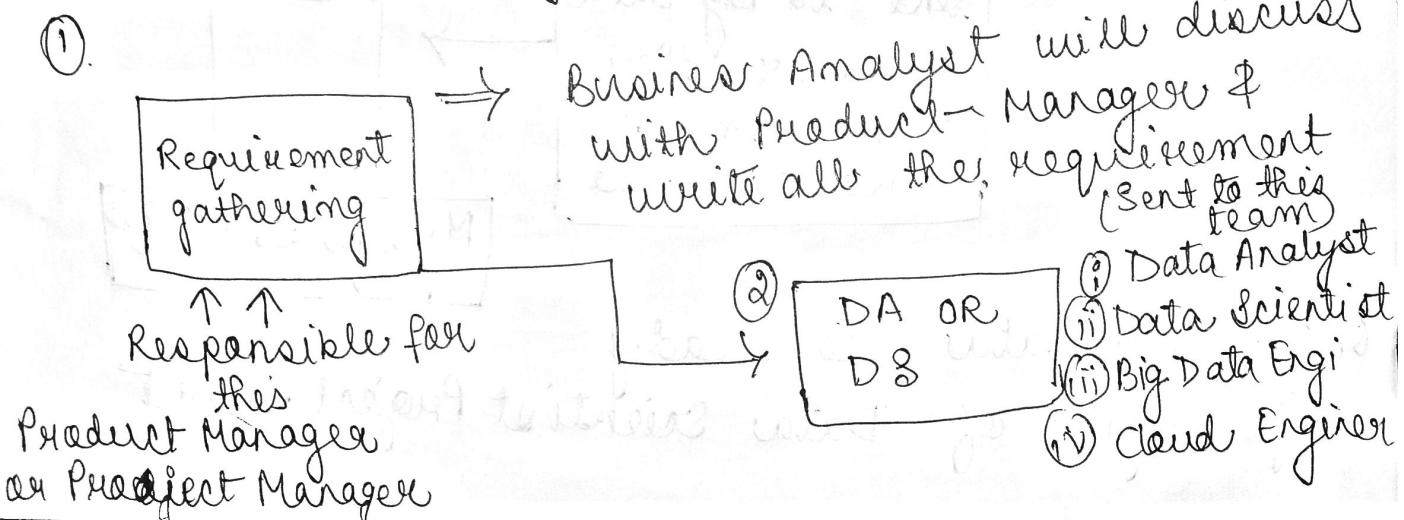
This task is basically few data analyst or data scientist

UseCase ② : Find the avg. size of the shark throughout the world?

Amazon Interview Question for DA

③ Amazon Big Billion Day Sale
asked in Intuit (product based) company
* which month should you select for Big Billion Day sale?

Statistics. [Life cycle of DATA SCIENCE PROJECT]



Service Based Company



e.g. T.C.S., HCL

giving service to
client.

Product Based Company



creating on their

own product

e.g. Apple, Google,
IBM

- ② DA or DS will tag
Product Manager and
BA will require here
to discuss to where
the data came from
with DA or DS.

Google

↓
Sales (department)

↓
Domain Expertise

↓
role: Product Manager

DATA CAN Get from

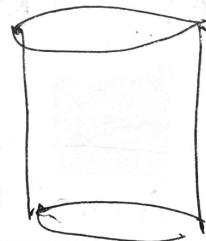
Internal
Database

3rd party
API's

Web
Scraping

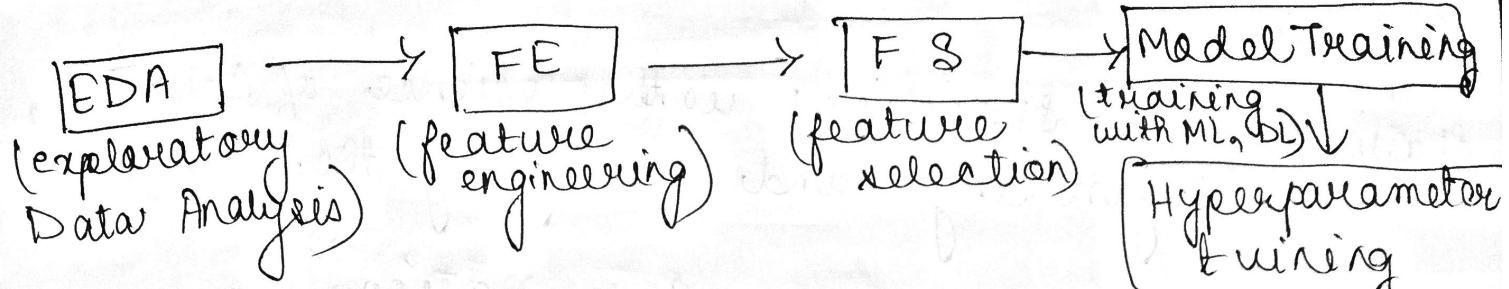
Combine

sent to Big Data
Engineer, they
will combine
and save to



MySQL or NoSQL

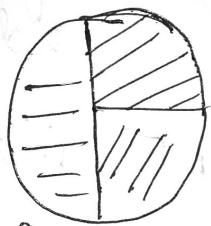
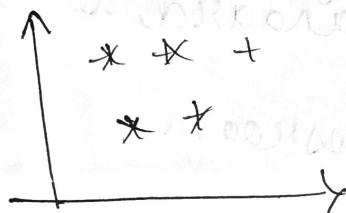
Once the data is ready
life cycle of Data Scientist Project start



* In all the steps Statistics will be used!!.

(Improve the performance of model)

Analysis of Data



Pie Chart

Summarisation of Data

Descriptive Statistics

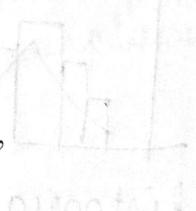
* Age = {12, 13, 14, 18, 20, 25}

Avg. age is a part of Descriptive Stats

measure of Central Tendency

Interview Question

How Statistics is used in Machine Learning!!!



Definition

Statistics :- Statistics is the science of collecting, organising and analysing the data.

Data :- "facts or pieces of information"

e.g. ① Ages of students in a classroom
 $\{24, 25, 32, 29, 18\}$

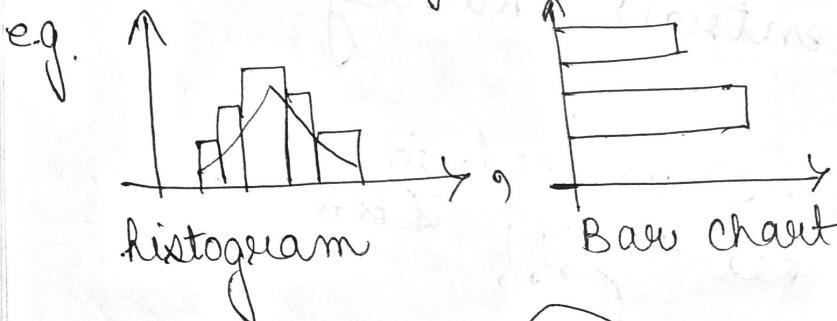
with the help of this data we can analyse this data by finding mean, median, mode, standard deviation

e.g. ② weight of students in a classroom.

Type OF STATISTICS

Descriptive Stats [EDA+FE]

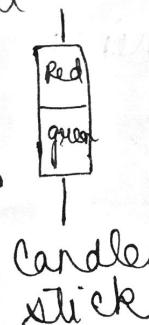
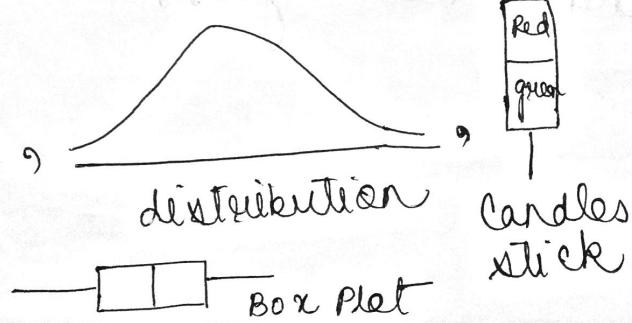
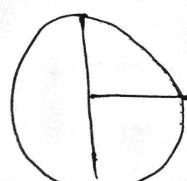
(1) It consists of organising & summarising the data



Inferential Stats

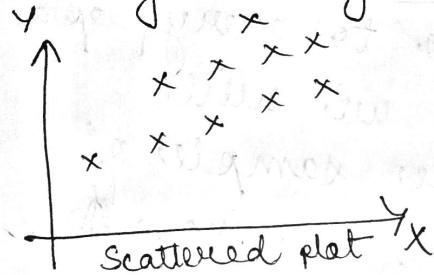
(2) It consists of collecting sample data and making conclusions about population data using some experiments.

* * making conclusions can be done by Hypothesis Testing



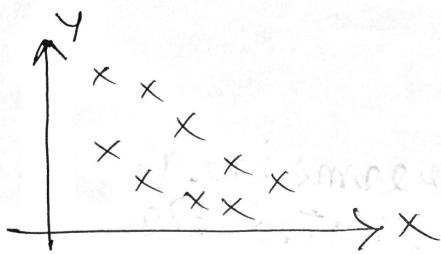
e.g. of Descriptive Stats e.g. of Inferential Stat

Let x & y features
↓
height weight



→ when $n \uparrow$, $y \uparrow$

→ when $n \downarrow$, $y \downarrow$



→ when $n \downarrow$, $y \uparrow$

→ when $n \uparrow$, $y \downarrow$

Let we have a

University of 500 people,
take a 60 people of

class A

↓
[Sample data] \rightarrow [Age]

↓
Avg. age of
entire university

Hypothesis Testing

c. I \Rightarrow Confidence
Interval

p-value

① z-test

② t-test

③ chisquare test

④ F test

we are taking

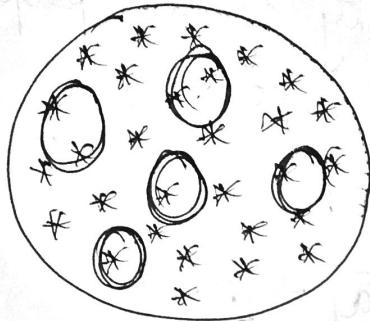
[Sample data]

& making conclusions

about [Population data]

through Hypothesis
testing.

Sample data Vs Population Data



Punjab State
(10 cro. popn)

Population Data

News will conduct some exit poll.

We can't go to every people to ask, so we will select some sample of size 1000 and ask from them.

Eg: Let's say there are 20 classrooms in a university and you have collected the ages of students in one classroom.

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

Weight h :-)

According to Descriptive Stats:- we get quest's like

- (i) What is the avg. age of students in the classroom
- (ii) Relationship bet'w age and weight.

According to Inferential Stats:- we get questions like

- (i) Are the avg. age of the students in the classroom less than the avg. age of the students in the university?

e.g. ② 1000 students in University

Class A \rightarrow 50 girls 50 boys

avg. marks 95% 92%

~~Answer~~ Ans:- girls performing well than boys.

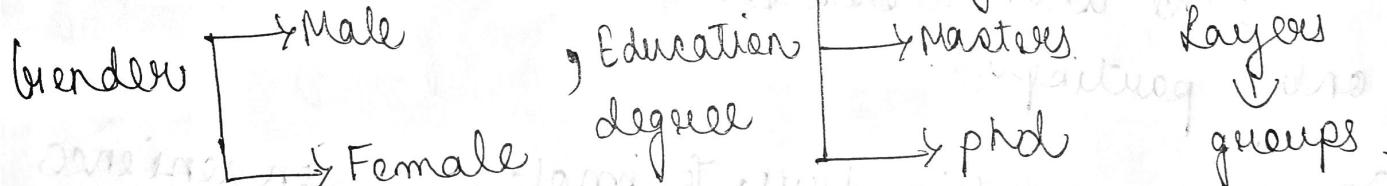
SAMPLING TECHNIQUES * Popn (N) sample (n)

(Choosing of sample techniques)

(1) Simple Random Sampling :- Every member of the population (N) has an equal chance of being selected for your sample (n)

c.g. exit poll, movie review, lottery

(2) Stratified Sampling :-



Blood group



In Exit Poll, we first apply stratified

sampling for neglecting the age < 18 ?

choosing the age > 18 [then for that pop we will apply random sampling]

Population [Exit Poll]

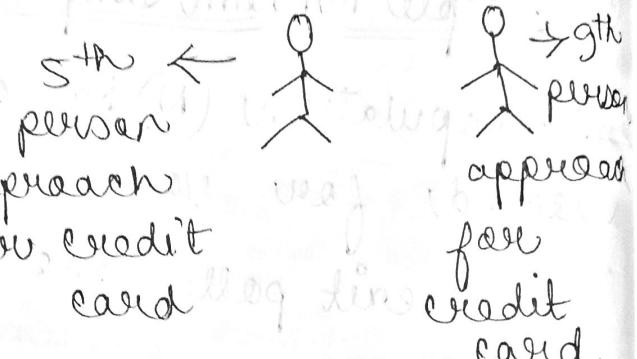
↓ Stratified Sampling
<18 year 18 year [Voting]

↓
Apply Random Sampling

③ Systematic Sampling

At {Airport}

Select every n^{th} individual
out of popn (N)



④ Convenience Sampling

Only those
who are interested in the survey will
only participate.

① Survey regarding New Technology \Rightarrow Convenience Sampling

② RBI \rightarrow conduct survey for **Women** \rightarrow **Married**
we will use ① Stratified sampling
then ② Random sampling.

③ Credit Card:- ① Stratified + Random Sampling
caller

① Variable : A variable is a property that can take any values
e.g. - age = 14, 25, 1000.

Variables

Ages = [24, 25, 26, 27] → Collection

Two different types of Variable

(1) Quantitative Variable :- Measured Numerically & Mathematical operators
e.g. Age, weight, rainfall (cm), temp^o, dist^r

(2) Qualitative Variables :- Categorical variables (Based on some characteristics)
e.g.; gender, Types of flowers, types of movies.

Quantitative



Discrete Variable

e.g.; whole number values

① No. of Bank Accounts person can have (1, 2, 3, 4, 5)

② No. of children in family

↓
Continuous Variable

e.g.; continuous
① Height, weight, ages, Rainfall, speed

- Assessment
- ① what kind of variable is Martial status
 - ② Yangtze River length
 - ③ Movie duration
 - ④ Pineapple
 - ⑤ IQ
 - ⑥ Gender

Ans:- ① Categorical Variable

② Continuous

③ Continuous

④ Discrete

⑤ Continuous

⑥ Categorical

PROGRAMS

Oldives - Oldives capitalised
intrants & exits

Great (ms) Alphabetic - View page 4
Oldives Designted - Oldives initialised
intrants & exits

Wellib - Page 4
Leverage Page 4

Initialised

Wellib Page 4
Leverage Page 4

Initialised

Wellib Page 4
Leverage Page 4

Initialised

Wellib Page 4
Leverage Page 4

Initialised