# Statistics
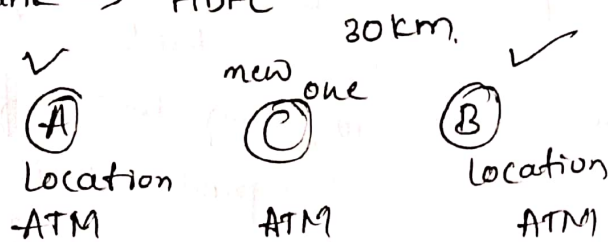
X statisticians ⟶ 5 years

→ Data analylist ✓
— Data scientist ✓
— Business analyst ✓
[ product manager ] — Domain expertise

TCS - servic

goyle - prod.
↓
creating
own
product

Usecase!- marking

Bank → HDFC

30 km.

✓        new one        ✓
Ⓐ        Ⓒ        Ⓑ
Location        Location
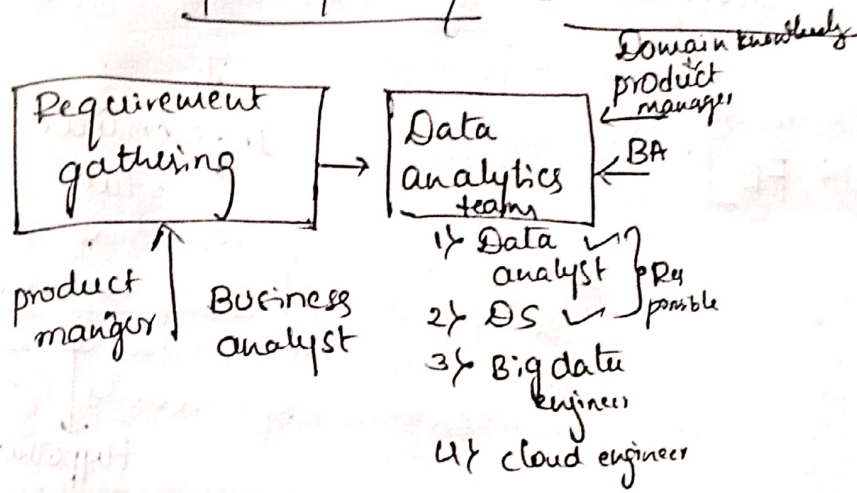-ATM        ATM        ATM

## Data engineering

→ Opening new ATM in C    How we can be come to conclusion we built ATM?

1) find the average size of the shark throught
the world ?    people
2) Amazon big billion day save { intuit }
which month should you select for the

## Statistics

1) Def^n
2) type of statistics
3) Life cycle of Data scientist
4) sample data (n) v/s population data (N)
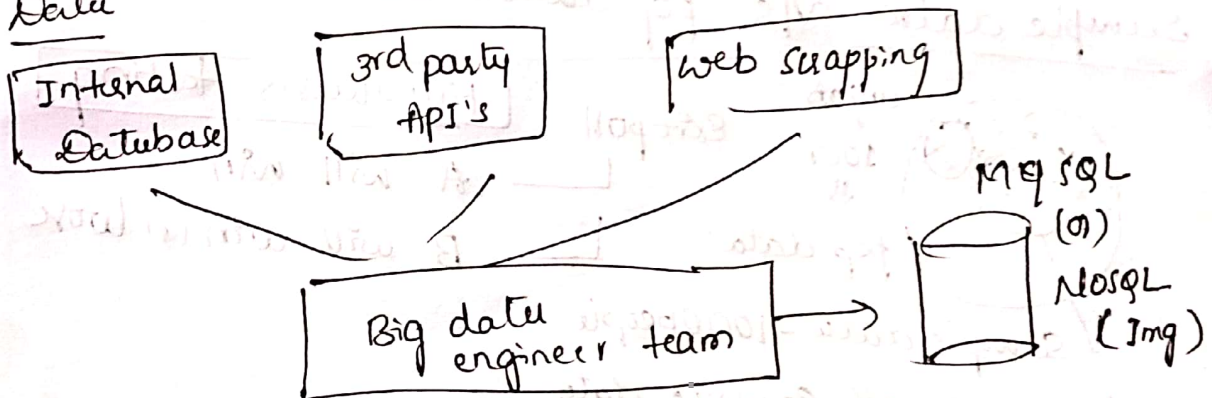5) sampling technique
6) variable types

# Life cycle of Data Science

Requirement gathering → Data analytics team ← BA ← product manager (Domain knowledge)

product manager | Business analyst

Data analytics team:
1) Data analyst ⌉
2) DS ⌉ PRg possible
3) Big data engineer
4) cloud engineer

## Statistics

It is The science of collecting, organizing & analyzing the data.

**Data** :- facts (or) pieces of information.

Eg :- Age of students - { 24, 25, 36 }

### Data

| Internal Database | 3rd party API's | web scrapping |

↓ → Big data engineer team → MySQL (or) NoSQL (Img)

## Data scientist :- project lifecycle

EDA → FE → FS (election) → Model training → Hyper parameter tuning

↑ Statistics

train ML (or) DL Alg

. Improve

- Analysis of data - visual

x̄ → summary data
|
Discrite statistics

Age = { 1, 2, 5, 8 } — Disciptive.
measure of central tendency

# Statistics

**Discriptive Stats** — Extensively used [EDA, FE]

**Inferential Stats**

1) It consist of orgnizing & summarizing using all plots
histogram


bar chart

pie | bell curve

candle stick | Box plot
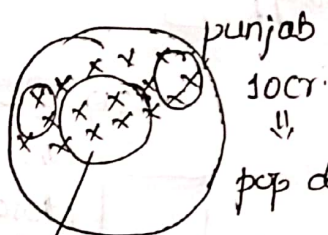
Conclusion
x? y?

→ It consist of collecting sample data & making conclusion about population data, ufing some experiments

→ Average calculation → Hypothesis test

| sample data | conclusion → | popl data |

Hypothesis testing.

(n)
## Sample data v/s population data (N)

punjab 10cr.
⇓
pop data

↳ Sample data — 1000 people

→ Average of all sample data

Exit poll — **Hypothesis testing**
↳ A will win
↳ B will win (a) lose

**Example:-** Let say there are 20 classroom in the university and you have collected the age of students in one classroom?

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

## Question

**Discriptive stats**
1) what is the average age of student in class?
2) Relationship blw age & gender?

## Inferential Stats:-

1) Are the average of the student in the classroom less than the average age of student in the university?
[ comparing - hypothesis test ]
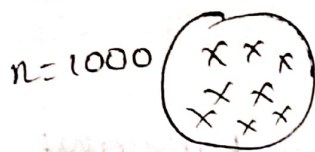
2) ∞ Average marks

    girls            50 boys      1000
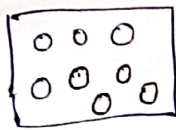    95 %.             92 %.      students

## Sampling Technique:-

1) simple random sampling:- Every member of the population (N) has an equal chance of being selected for your sample (n).

$n = 1000$    (x x x / x x x / x x x)

Exit poll          } Random sampling
genral servey

eg :- ☐ (o o o / o o o o)    select the marble, each marble have equal chance of selecting.

2) stratified sampling:- [ strata → layers → clustin → group ]

① gender [ F / M ]    ② education [ High school / master ] degree

electing

③ [ Exit poll ]    vote →18 → Random sampling
stratified    not <18 vote →    ④ blood group.
sampling.    └→ R.

3) systematic sampling:-
   { credit card }    → select every $n^{th}$ individual out of population (N).

5th           9th
person ☺ ☺    person

→ This is slm way of selecting. ($n^{th}$)

4) convience sampling: only those who are interested in survey will only participate.

Ex: DS survey → general AI. who interested

→ ineaion job for specific role - those who are really interested those only fill.

1) survey regarding new technology.
          ↳ convience sampling.

RBI → survey → women. ⟹ stratified + Random
                ↳ Married

2) Credit card call
         ↳ stratified + Random
           salary..

Variable :- It is a property that can take any value.
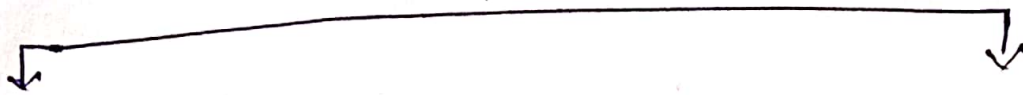
variables

Ex. Age - 14
    A = 24

Aages = [1, 2, 3]

2 - types of variable

1) Quantitative variable - measure Numerically.
                    { maths opuation }
      Ex: Age. weight, temp, distance

2) Qualitative variables - [ categorical variable ]
    based on some characterstics they grp together
Ex: gender, flower, car, movies, dept.

# Quatitative

**Discrete**
Variable [many
              no of
              variable]
              categorical)

Ex. whole no
- No of bank acc
  { 1, 2, 3, 4 }
-> limit to no
   again repeat
   w.r.t data points

**Continous vaid**
Ex continous.

Ex. Height, age, speed,
    Rainfall
-> decimal value.

1) Marital status ?  — categorical [ no of maried people
                                        discrete]
2) ganga river length ? continous
3) movie duration ? continous
4) pincode ? [descrete - continous] not categorical bcz
5) IQ ? continous (discrete) more no so not able
6) gender ? categorical          analyze instead
7) country ? - discrete          descrete - Feature
                                        engineering