# Statistics

Can work as
- → Data Analyst
- → Business Analyst
- → Data Scientist
- → Product Manager

Usecase:

Asked in Product Based Company ?

1) HDFC

Check whether to open this ATM or not ?

Loc 1
A
↑
ATM

C
↑
New ATM

Loc 2
B
↑
ATM

This ques can be solved by
① Data Analyst
② Data Scientist

←——— 30 km ———→

or Blue Whale ← (Amazon)

2) Find the Avg. size of the [shark] throughout the world.

← (Intuit)

3) Amazon Big Billion Day sale. Sortily → which month should you select ?

Topics to cover in statistics :- { Life cycle of Data Science }

Domain knowledge    Project

[Data science] → Data Analytics Team

Requirement Gatering → 
- Data Analyst ← Product manager
- Data Scientist ← BA (Business Analyst)

① Data Analyst
② Data Scientist
③ Big Data Engineers
④ Cloud Engineers

↑ ↑
Product manager    Business Analyst
↑
Project Manager

Types of company
Service - TCS, Infosys,
Product Based
(FB), {Google}, {Apple}
YT, G Pay
Google Ads, Gmail

[Sales]
↑
Domain expertise
↓
Product manager

# From Where Can you get the Data

```
┌──────────┐      ┌──────────┐      ┌──────────┐
│ Internal │      │ 3rd Party│      │ Web      │
│ Database │      │  APP's   │      │ Scraping │
└──────────┘      └──────────┘      └──────────┘
        ╲              │              ╱
         ╲             │             ╱
          ╲            ▼            ╱
          ┌──────────────────────┐
          │  Big Data Eng.       │        ┌────┐    mySQL
          │  ___Team___          │  ───▶  │    │     or
          └──────────────────────┘        └────┘    NOSQL
```

## Life Cycle of DS Project

Exploratory Data Analysis    Feature Engg.    Feature Selection

```
┌──────┐      ┌──────┐      ┌──────┐      ┌──────────┐      ┌──────────┐
│ EDA  │ ──▶  │  FE  │ ──▶  │  PS  │ ──▶  │  model   │ ──▶  │ Hyperparam│
└──────┘      └──────┘      └──────┘      │ Training │      │  Tunty    │
   ⇑                                      └──────────┘      └──────────┘
                                               │                 │
Statistics                               Training with      Improve the
                                         ML Algo.           Performance
        ⇑   ↑                                or                 of
     Analysis of                          D.L Algo.           model
        Data
```

Where Do we use statistics?

```
  │
  │  x
  │ x x  x
  │  x  x   x
  │ x x x x    ⟹      Discriptive stats
  └──────────            ⇑
                    Summarizing the Data
```

Pie Chart

```
   ╭─────╮
  │   │   │
  │───│───│    ⟹    Distributive Stats
  │   │   │
   ╰─────╯
```

Age = $\{12, 13, 14, 18, 20, 25\}$ ⇒ Avg. Age ?

⇓ → Measures of
central
Part of Tendency
Discriptive
stats.

Statistics Def^n : → Statistics is the science of
collecting, organizing and analysing the Data.

Data :→ "Facts or pieces of information"

eg: (i) Ages of Students in classroom.

$\{24, 25, 32, 29, 28\}$ ⇒ Mean, Median,
Mode, Standard
Deviation,

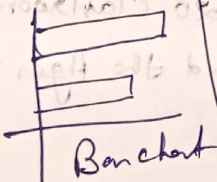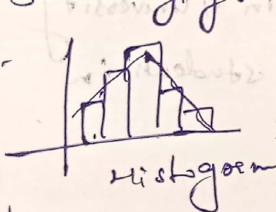(ii) Weights of Students in classroom.

## Type of Statistics

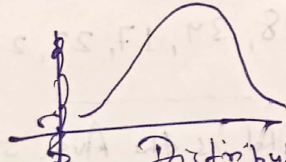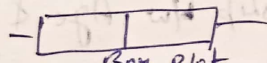Descriptive stats

(i) It consists of organizing &
summerizing the data.

eg -



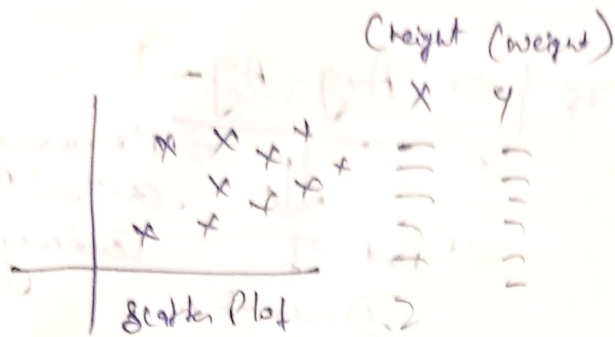Histogram , Bar chart



Pie chart

Candle stick

Box Plot

Inferential stats

(*) It consists of Collecting
Sample data & making
Conclusions about Population
Data using some experiments.

→ Hypothesis
Testing.

University

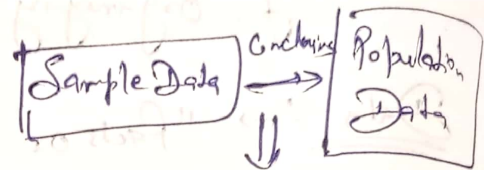Class A → 60 People

⇓

Sample data ⇒ Age ⇒ Age
of entire university

(height (weight)

X        Y

Scatter Plot

| X ↑ | Y ↑ |
| X ↓ | Y ↓ |

| X ↑ Y ↓ |
| X ↓ Y ↓ |

Hypothesis Testing

C·I ⟹ Confidence Interval

P- value

① Z- Test
② t - test
③ chi square test
④ F Test (Anova Test)

| Sample Data | → Concluding → | Population Data |

⇓

Hypothesis
Testing }

---

# Sample Data Vs Population Data.

← Punjab

Exit Poll
↑
Party A will win }
Party B will lose }

If opp. happens
means hypothesis
is gone wrong.

| 10 cr |
⇓
Population
Data

Sample
size > 1000

Eg: Q. Let's say there are 20 classrooms in a university
and you have collected the Ages of students in
1 classrooms.

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

weight = { _____ }

Descriptive stats :→ What is the Avg. of students in classroom?
Relationship b/w Age & weight?

**Inferential stats** :→ Are the avg. age of the students in the classroom is less than the avg. age of students in the university

→ greater than

↓ equal to

Eg :- **University has**
1000 Students

<u>Sample Data</u>

Class A → 50 girls    50 boys

↓ Avg. marks   ↓ Avg. marks

95%    92%

& Can we Conclude that Girl has done better than Boys in the entire university?

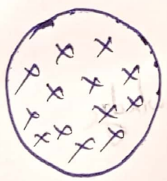→ Yes for this we have to Apply
    Hypothesis Testing.

⇒ **Different Sampling Techniques** :→   | Population denoted by (N) |   | Sample denoted by (n) |

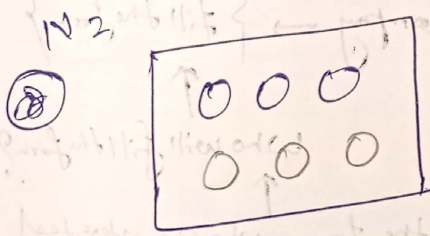① <u>Simple Random Sampling (SRS)</u> :→ Every member of the population (N) has an <u>equal chance</u> of being selected for your sample. we can use simple

Eg :-

Crit :→ Random Sampling
Poll
General Survey,
Movie Reviews,
Lottery.

N 2,

② $n = 3$

Every marble has equal chance of getting selected.

② <u>Stratified Sampling</u> :→ | Strata → Layers → Clusters → Groups |

(i) Gender → Male / Female

(ii) Education Degree → High school / Masters / PHD

(iv) Blood groups → O+ / A+ / B+

(iv) Exit Poll → <18 / >18 ✓ ← Apply Random Sampling

③ Systematic Sampling :→ { AIRPORT }

{ Credit Card }

Approach
every 5th
Person
for credit
Card

Approach
every 9th
Person for . .
~~Apply credit~~
Card

Here we select every nth individual out of Population (N).

Eg: every 5th Person, or every 10th Person etc.

they are following some systematic ways,

④ Convinience Sampling :→ Only those, who are interested in
the Survey will only participate

Eg:

ⓘ { Data Science :- General AI Survey }
Survey

↑

Whoever Participate in this Survey
Should be interisted or have knowledge
about Data Science

& ⓘⓘ  iNeuron Job for specific Company → { fill the form }

↑

Who will fill the form?

↑

Those will fill the form who are intersted
for that Particular Job.

8. What Sampling Can we use in the following Situation?

① Survey regarding New Technology → Convinience Sampling — Bcz (only those who are interested technology will participate)

② RBI Survey by → Women → Stratified Sampling + Random Sampling
(Since women takes care of entire house)
   └→ Married Women

③ Credit Card Call :→ Stratified + Random ( Bcz, mostly they call Salaried people & then Random )

⑨ → Variable :→ It is a property that con take any value.

Eg: age = 14          Variables          √ list, collection
    age = 25          Ages = {24, 25, 26, 27, 28, 29}
    age = 100

⇒ Two Different types of Variables :→

① Quantitative Variable :→ Measured Numerically, { Mathematical operations }

   Eg :- Age, weight, height, rainfall, temp, distance

② Qualitative Variables   (Catogrical Variables)
                               ↓
                          ( Bcz of Based on some Characterics they are grouped together )

   Eg :- Gender, Types of Flowers, Types of movies etc.

# Quantitative Variables (Continuous Variable)

## Discrete Variable → Fixed

Eg :- [Whole no.] (decimal nos. not allowed)

i) No. of Bank a/c.
$\{1, 2, 3, 4, 5\}$ [2.5] X

ii) No. of children
Should have less no. of categories.

## Continuous Variable → Decimal Value

Eg :- Values will be [Continuous]
→ Any value;
Eg :- Height, Weight, Ages,
Rainfall, Speed etc.
Can have whole nos. also,

## → Assessment

① What kind of Variable is

→ married
→ Not married

Marital status ? → Categorical Var

Ganga River Length ? → Continuous

Movie Duration ? → Continuous

Pincode ? → Discrete
↳ 105, 75, 90.5

? & ? → Continuous Discrete
Since it doesn't have decimal values;

Gender ? → Categorical

No. of People married ? → Discrete

## Pincode

160099
560098
560097
} unique

↓
PANCARD → Categorical
Variables

Since # is
having Alphanumeric
values,

## Final

It is many ?

So we can't take
it as Categorical.

[Categorical]
⇓
[Feature Engg.]