# Lecture 1

## Statistics

For ML, DL, NLP, Vision, Data Analyst

## Definition

It is the branch of science which involves collecting, analysing, data in large quantities so that we can come up with various use cases, exploring and visualising and coming out with meaning information and conclusions

Topics which we will see:

Descriptive statistics
Inferential statistics
Population
Sample
Sampling techniques
Measure of central tendency
Measure of dispersion or Variance
Probability
Permutation and combination

# Statistic

## Descriptive

① Analyze, Explore

Visualizing → Techniques

↓

> To understand the data

Eg; Histograms, bar, pie, scatterplot

## Inferential

① Population of data

↓

Sample of data

↓

Experiment

↓

> Conclusion Find Out

Eg; Hypothesis Testing

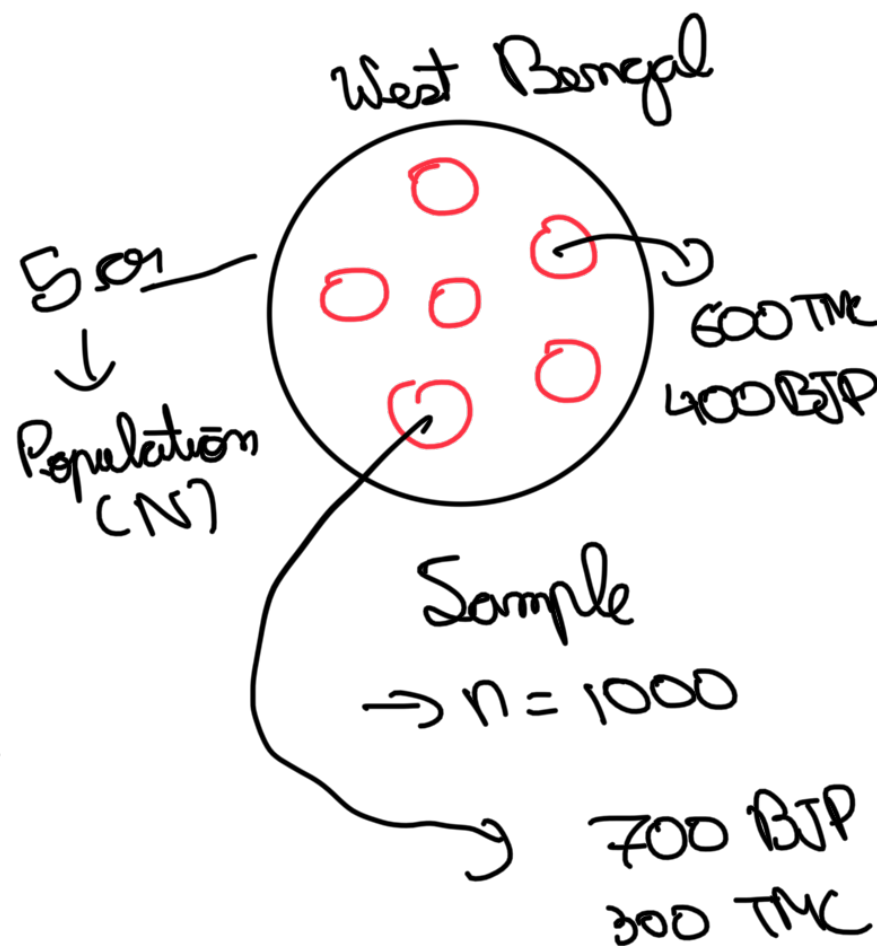P - Test,
T - Test,
F - Test,
Z - Test

Null Hypothesis

Alternative

# Decision Making

## Population (N)

From these samples we can say

{ Max seats will go to TMC, and second highest will be BJP }

West Bengal

500 → Population (N)

600 TMC
400 BJP

Sample
→ n = 1000

700 BJP
300 TMC

---

Different sampling methods

1. **Random sampling-** randomly getting selected
**Disadvantage :-**
a) overlapping- some cases are being repeated
b) For specific use case -

2. **Stratified sampling-**population is divided into Stratus( male and female)
Male for exit pole
Household for females

So we are dividing them into subgroups.

3. ***Systematic sampling-*** doing sampling systematically, that is if we take the second person then 4th then 6th then 8th. We are following a system

It might lead to bias which means more than 1 type in a sample

4. ***Clustering sample-*** Clusters of groups

Every group will give different answers

We will take groups instead of individuals

Targeting different clusters of customers( based on expenses we will target different sets of rich and poor customers)

# Measure of Central Tendency

① Mean  ② Median  ③ Mode

$$\bar{x} = \frac{\Sigma n}{c}$$

$\{1, 2, 3, 4, 5\}$
mean = 3

If $\{1, 2, 3, 4, 5, 100\}$
mean = 19.6

100 has become outlier

When we have outliers then we should never take mean

Q If we should not use mean what should we use?

⇒ We can use median.

$$\{1, 2, \boxed{3, 4}, 5, 100\}$$

$$\text{Median} = 3.5$$

$$\{1, 3, 100\}$$

$$\text{Median} = 3$$

---

Median does not get impacted with outliers ✳

---

In median we need to always sort the elements. ✖

---

③ **Mode** – **Most frequently occurring element**

334 → Mode 3

1333444 → 2 modes
$$\boxed{3 \text{ and } 4}$$

In latest python the first by default which is present more

times that it is the **mode**.

---

_Random variable_ - whose value depends on the outcome of a random phenomena

Tossing a coin :.

$\Rightarrow$ we can get $\{H, T\}$ which is a random phenomena

a) Categorical Variables
   Qualitative Variables

b) Continuous Variables
   Quantitative Variables

a) $\Rightarrow$

| Gender | Weekdays |
|--------|----------|
| └ M    | Sun      |
| └→ F   | Mon      |
|        | ⋮        |
|        | Sat      |

a) ⟿ Nominal → Rank not important
            eg F and M, etc
   ↝ Ordinal $\{$→ Ranking needs to
              be considered

┌─────────────────────────┐
│ For IT Professional      │
└─────────────────────────┘

Sunday , Sat , Fri , . . . . . .

Randenig

| Customer Ratings |
|---|

5, 4, 3, 2, 1

---

b) Continuous Variable

eq; Height = { 170.3, 171.4 , 180.;
2 }

↱ Discrete Quantative Variables
↳ Continuous Quantitative Variables

| Temperature | — Continuous Quantative Variable |
|---|---|

| Pincode | — Nominal Categorical |
|---|---|

| Age | — Descriptive Quantative |
|---|---|

---

Variable

1) Quantative Variables
2) Qualatative Variables
3) Random Variables

---

When there is NaN value we cannot replace it with mean or median coz we need numerical value in that case so we use mode

---

For mean → Remove Outliers

↘ Or else use median

---

① Independent Samples
② Dependent Samples

---

| Σ Temp | Rainfall | $CO_2$ | $NO_2$ | Humidity |
|--------|----------|--------|--------|----------|

Independent features

With this we calculate

Output

AQI

↳ Dependent feature