# STATISTICS

## Introduction to Stats:

Statistics is the branch of applied mathematics that involves the collection, description, analysis and inference of conclusions from quantiative data.

Statistics is the Science of the data to Collect, Organize and analysis to made Conclusions or decision Making.

| Descriptive Stats | Inferential Stats |
|---|---|
| It focus on Collecting, Organizing and Summarizing the data. | It focus on to made Conclusion. |
| | It is the technique where in we Used the data that we have Measured to form Conclusion. |

# Descriptive Stats:

↓

1. Measure of Central tendency.

2. Measure of Dispersion.

## Summarizing the data,

Histograms, pdf, cdf, Probability,

Permutation, Mean, median, mode, Variance, Standard

deviation.

## Distributions,

1. Gaussian (or) Normal distribution.

2. Log Normal Distribution.

3. Binomial Distribution.

4. Bernauli's distribution.

5. Poisson distribution.

6. Pareto distribution.

7. Transformation and Standardization.

8. Q-Q Plot.

# Inferential Stats:

$$\downarrow$$

(F-Test)

Z test, T test, ANOVA Test, CHISQUARE,

Hypothesis test (p values), Confidence Interval,

Z-table, t-table.

## What is Data?

Facts or pieces of information that can be measured.

Eg: IQ of students = $\{90, 80, 50, 100, 120\}$  ← Data

Student's marks = $[90, 80, 85, 70, 60, 30, 55, 85, 93]$

↳ one particular math classroom marks

The common question from above marks in the point of descriptive stats are,

1. What is the average students marks?
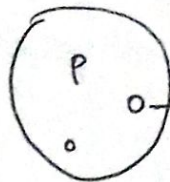2. How many students are get marks above 50?

etc.

The common question from above marks in the point of inferential stats,

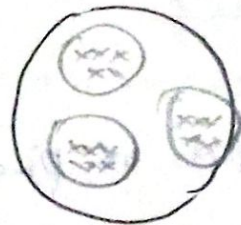1. Are the marks of the students of this classroom similar to the number of maths classroom in the college?

whole math classroom

Particular one math Classroom

$P$ $o$ $o$ → Sample

Population (N) - Entire population among the data.

Sample (n) - Some Sample from entire population of the data.

<u>Some Sampling Technique</u>:

<u>Simple Random Sampling</u>:

Every number of the population (N) has an equal chance of being Selected for your Sample (n).

It focus on random Sampling.

## Stratified Sampling:

Where the population $(N)$ is split into non-Overlapping groups. (strata)

Ex: We want to collect and sample about Job Professions. We need to sample get from only Job Profession with respective field doctor, engineer, IAS etc. [Not applicable for collect Sample from whole population]

Ex 2: Gender $\left[\begin{array}{l} \rightarrow \text{male} \\ \rightarrow \text{female} \end{array}\right.$

Non-Overlapping groups - One group independent from another group.

## Systematic Sampling:

$(N) \rightarrow n^{th}$ individual

Eg: Mall - Survey about Covid for every 10th person Out from the mall.

Systematic Sampling depends on the nature or Scenario.

## Convenience Sampling:

Samples are collected from population if only for convenience.

Example:  DATA SCIENCE
                        ↑
          only get Samples from DATA SCIENCE
                in Convenient way.

## Examples for Sampling:

VOTE POLL SAMPLE - Random Sampling used.

RBI Women Survey - Convenience Sampling.

DRUG TESTED      - Based on Condition and
                   nature.

## Variables:

A Variable is a property that can take on any value. [STATS DEFINITION]

A Variable is the Container to store the Value. [PROGRAMMING DEFINITION]

Both are
Same

Eg:

Height = {78, 65, 60, 40, 50}

Weight = {65, 80, 100, 120, 75}

Two kinds of Variables:

→ Quantitative Variable.

→ Qualitative Variable.

Quantitative Variable:

It can be measured and numerically type and supports Operation for Add, mul etc..

Quantitative Variable
                /            \
          discrete          Continuous
    Eg: Whole number [0,9,10]    Eg: height = [172.5, 162.5...]
    No of Bank accant -5         weight = [100kg, 99.5]
    [Cannot be split]            [Can split]

## Qualitative Variable / Categorical Variable:

It is based on the data other from Quantitative like strings.

Gender
$\begin{bmatrix} \to f \\ \to m \end{bmatrix}$
{ Based on Some Characteristics we Can derive Categorical Variable }

Example:

Gender, email etc...

## Variable Measurement Scales:

4 types of measured Variable.

→ Nominal [categorical data]

→ Ordinal [order matters, does not require Values]

→ Interval [order matters, Value matters, natural 0 are not Present]

→ Ratio [order matters, Value Matters, Natural 0 present]

Nominal ~ Male.

→ This ordered data is known as Ordinal.

Ordinal ~ (RANK) marks

| RANK | marks |
|------|-------|
| 1    | 95    |
| 2    | 80    |
| 3    | 70    |

Interval - [70-80] [80-100] but 0 not
present
50°        60°

0 → does not respect.

Ratio - It possesses all the properties of interval

data and allows for meaningful ratios

between values.

height : [180, 180, 170]

↓    ↓    ↓
A    B    C

The ratio of person B height to person C height

$= \dfrac{180}{170} = 18:17 \ (60) \approx 1.05883$

## FREQUENCY DISTRIBUTION :

→ It is the distribution to return and

Visualize the count through various graphs

like bar graph, Histogram.

→ Let consider the simple example,

Number of population in countries like.

## Sample dataset:

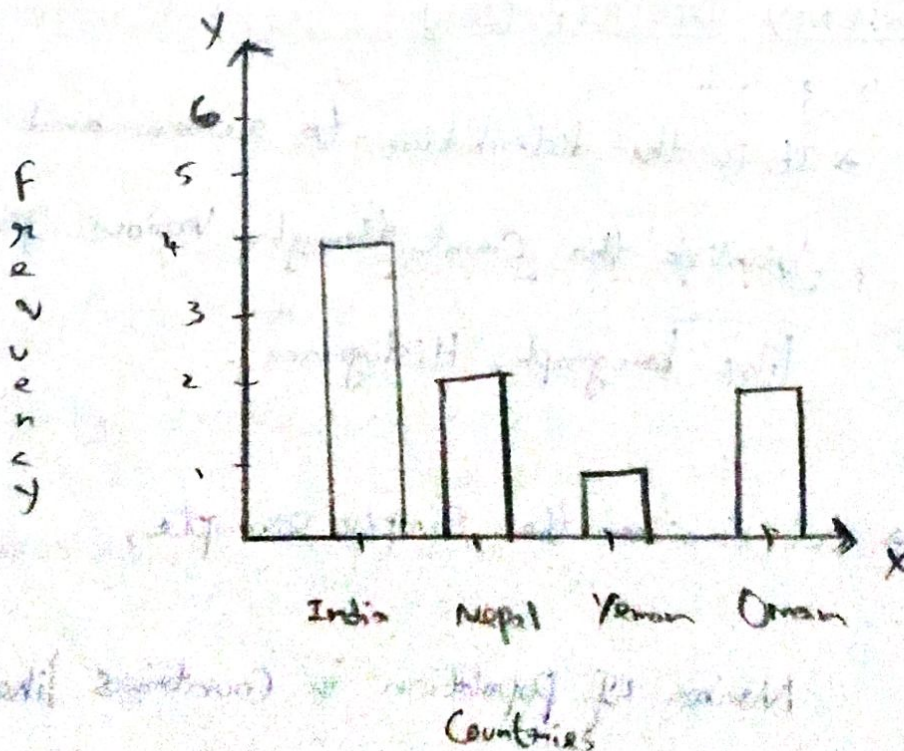Country = [India, nepal, nepal, india, india, india, Yeman, Oman, Oman].

| Countries | Frequency | Cumulative Frequency |
|-----------|-----------|----------------------|
| India | 4 | 4 |
| Nepal | 2 | 6 |
| Yeman | 1 | 7 |
| Oman | 2 | (9) |

↑ Total Count

Visualize through,

(Summarizing the data)

## BAR GRAPH:

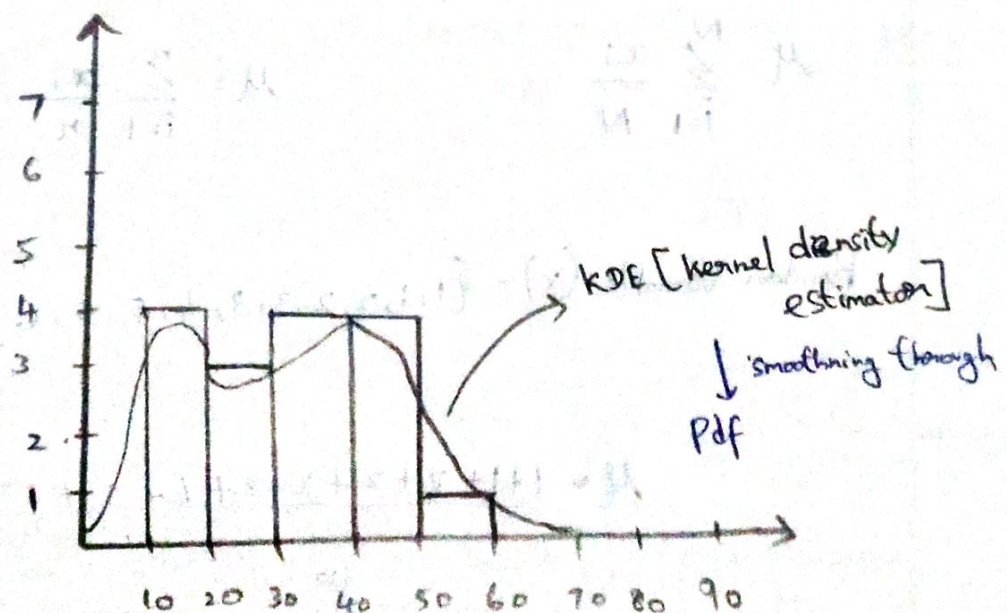Bar graph only represents the discrete data (or)
Categorical data.

Suppose use Continuous to use histogram,

## HISTOGRAM :

Sample data ,

$$age = \{10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41,$$
$$42, 43, 50, 51\}$$

Default histogram bin size = 10 [able to change]



KDE [kernel density estimator]
↓ smoothning through
Pdf

Pdf is nothing but smoothening of histogram.

Pdf [Probability density function]