# Stats

## Day-1

* Introduction to stats :-

1. Basic to Advance { Data Scientist, Data Analyst, Business Intelligent Tools }

→ First 2 days Basics

1. Descriptive Stats
2. Inferential stats

1. Descriptive stats
{ i. Measure of Central Tendency
  ii. Measure of Dispersion }

Any thing relating to Summarizing data.
→ Histograms, pdf, cdf, Probability, Permutation, mean, median, mode, variance, Standard deviation,

1) Gaussian Distribution
2) LogNormal Distribution
3) Binomial Distribution
4) Bernoulli's Distribution
5) Parito Distribution { Power law Dist }
6) Standard Normal distribution
7) Transformation and Standardisation
8) Q-Q plot

2. Inferential stats :-
Z-test
T-test
Anova
Chisquare
Hypothesis -Testing

Z, table, T table.

Now,

**Q. What is Startistice ?**

→ Startistics is the Science of Collecting organizing and analyzing data.

{ Better Decision making

Definition of data ?

→ Fact or pieces of informatra that can be measures.

Eg:- i) The IQ of a class a student
{ 98, 97, 60, 55, 75, 65}

ii) Age of Students of a Class
{30, 25, 24, 23, 27, 28} → Data

**✳ Types of Statistics :-**

i) Descriptive Stat.
→ It Consist of Organizing and Summarizing data.

ii) Inferntial Stats.
→ It is a Technique where we use the data that we measured to form Conclusion.

Eg:- Classroom of maths Students (20 student)
marks of the 1st Sem

84, 86, 78, 72, 75, 65, 80, 81, 92, 95, 96, 97. - - - - - - .

Eg:- Descriptive Stast:-
What is the average marks of the Students in the class.
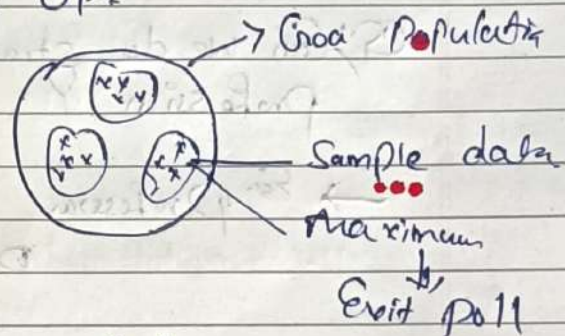
Eg:- Inferential Stat:-

→ Are the marks of the Students of this Classroom similarly to the age of the maths Classroom in the cBte college. ?

Population and Sample :-

Election → Goa, UP.

Evit poll



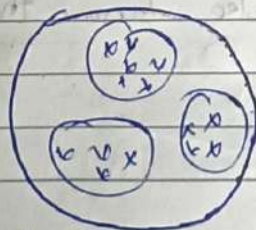Goa Population

Sample data

Maximum

Evit poll

Population (N)        Sample (n)

* Sampling Techniques :-

i) Simple Random Sampling :- ←



When Performing simple random Sampling Every no. of population (N) has an Equal chance of being Selected. for your sample N.

ii) **Stratified Sampling :-** is a technique rather the population is split into Non- Overlapping groups (Strata)
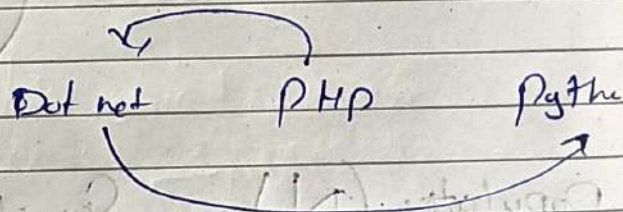
Eg:-

Gender 
- → Male
- → Female

Survey.

Ex :- Age-group

(0-10)    (10-20)    (20-40)    (40-100)

→ **Interview Question :-**

● Q) Can we do stratified Sampling based on Profession ?

••• → Eg:- Profession

Dot net        PHP        Python

⇒ Some of the Stages we can do Stratified Sampling

→ But, applying stratified Sum other Condition we can make Sure that the Sampling Satisfies.

iii) Systematic Sampling :-

(N) ⟶ nth individual

Eg :-
mall → Survey (covid)

↳ 8th person ⟶ Survey
↳ 10th person → Survey

( Indipendantly Survey )

iv.) Conveince Sampling :-

Eg :- ↳ Survey
        ↳ Data science ← Only these
                            people
{ Basically interest / Expert in DS }

→ Survey related to Specific topic.
  → In this case Data Science.

Eg :- Exit poll
       { Random Sampling using }

Eg :- RBI → House hold Survey
                      #
            Survey → women
{ Stratified Sampling / Conveince Sampling }
Eg :- Drug → Tested ⟹

* **Variables :-** A variable is a property that can take on any value

Eg:-

Height    { 182, 177, 168, 180, 175 }
Weight    { 78, 99, 100, 60, 50 }

Two kinds of variable :-
1). Quantitative Variable
11) Qualitative Variable / Categorical Variable

⇒ 1) Quantitative Variable :-

Measured Numerically { Add, Substract, mul, div } v

Eg:- Age
weight
Height

11) Qualitative / Categorical variable :-

Eg:- Gender $\begin{bmatrix} M \\ F \end{bmatrix}$

→ Based on Some characteristic we can divide Categorical variable

Eg:- IQ

0-10            10-50            90-100
↓               ↓                ↓
Less IQ         Medium IQ        high IQ
                                 Good IQ

Eg: Blood group        Tshirt

A+           XS

A-           S

O+           M

O+           L

Etc           XL

          XXL, etc

⇒

Quantitative

Descrete variable        Continuous variable.

Eg:- Whole number        ⇓

Eg:- No of Bank Acc      Eg:- Height

→ 2, 3, 4, 5, 6,7        172.5, 162.2, 172.3, Et

Eg:.

→ +b .of childrens in a family    weight.. 100, 99.5,

Eg:. 3, 4, 5, 2, 3,          .77.5, Etc.

Amount of Rainfall:-

1.1, 2.3, 1.35, Etc.

Sample Questions:-

1. What kind of Variable Grad is? ⇒ Categorical

2. " " " " marital status? ⇒ "

3. " " " " Population of State is? ⇒ Discre

4. " " " " River length? ⇒ Continues

5. " " " " Song length? ⇒ Continuous Continu

6. " " " blood pressure? ⇒ Discrete Continu

7. " " " PIN Code? Discrete

**\* Variable Measurement Scales :-**

4. types of measure of variable

1. Nominal → { Categorical data } → Classes $^{Eg:- Color, Gender}$
2. Ordinal → Order of the data matters, Value does not
3. Interval → Order matters, value matters, $^{natural zero not present}$
4. Ratio.

Eg:- 2. ordinal :-

| Students (marks) | Rank |
|---|---|
| 100 | 1 |
| 96 | 2 |
| 52 | 4 |
| 85 | 3 |
| 44 | 5 |

→ ordinal data

3. Interval :- Order matters, value also matters, natural zero not present

Eg:- Temperature
→ Fahrenheit

| 70-80 | 80-90 | 90-100 | 100-110 |
|---|---|---|---|

4. Ratio :-

Absolute zero points, Equal intervals, Ratio are meaningful, Quantitative

eg:- length (measure in meters, cm, etc)
Weight ( " in kg, grams etc)
Time ( " " second, min, hrs, etc)
Income ( " " currency, dollars, etc)
Age ( " " years, months, days, etc)

## * Frequency Distribution :-

Sample datast :- Rose, lilly, Sunflower, Rose, lilly,
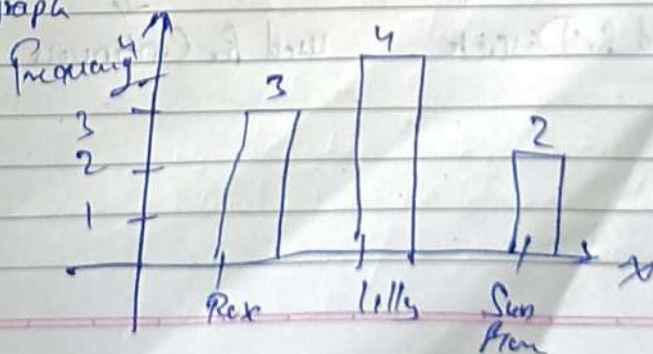Sunflowers, Rose, lilly, lilly

| Flower | Frequency |
|--------|-----------|
| Rose | 3 |
| lilly | 4 |
| Sun flower | 2 |

→ Frequency distribution table

→ Cumulative Frequency

| Flower | Frequency | Cumulative Frequency (CF) |
|--------|-----------|---------------------------|
| Rose | 3 | 3 |
| lilly | 4 | 7 (4+3) |
| Sunflower | 2 | 9 (7+2) |

eg:- in graphs

i) Bar graph
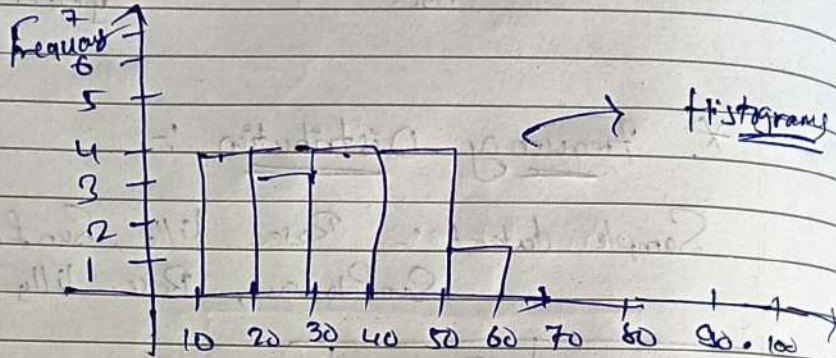
Eg:- II) Histogram :- Continuous :-

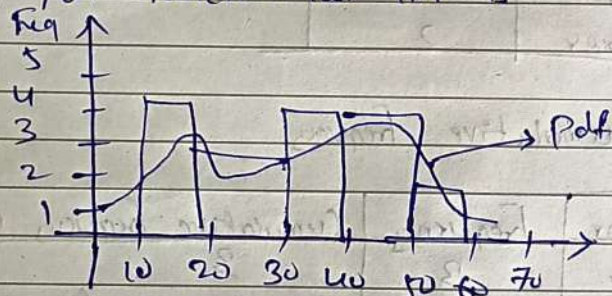Age = { 10, 12, 14, 18, 24, 26, 30, 35, 38, 37, 40,
          41, 42, 43, 50, 57 }

Bin = 10 ┘
   ↓
(Grouping)

Frequency

7
6
5
4
3
2
1

10 20 30 40 50 60 70 80 90 100

→ Histogram

Pdf :- Smoothning of histogram
My pdf function look like ↓

Freq

5
4
3
2
1

10 20 30 40 50 60 70

→ Pdf

KDE { Kernal Density Estimator }
PDF { Probability density function }

* BAR  V/s  Histogram
    ↓              ↓
Used for Discrete   Used for Continuous