

Introduction to statistics

Data & Statistics is the science of collecting, organizing and analyzing data

Data :- "facts or pieces of information"

Eg:- Height of students in classroom

175cm, 180cm, 190cm --

Types of statistics

1) Descriptive statistics

→ It consist of organizing and summarizing data

2) Inferential statistics

→ It consist of using data you have measured to form conclusion

- Measure of Central Tendency

[Mean, Median, Mode]

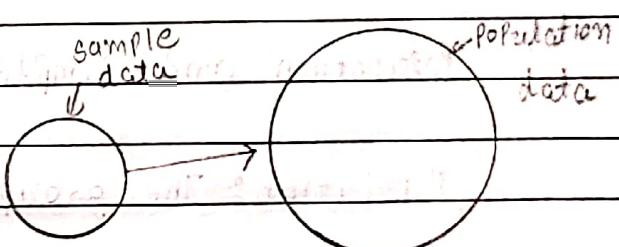
- Measure of Dispersion

[Variance, Std]

- Different type of distribution

of data

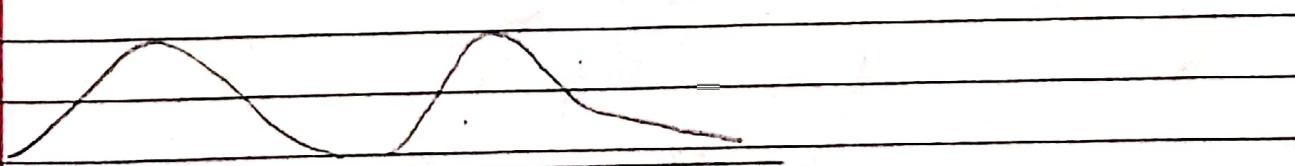
Eg:- Histogram, P.d.f, P.m.f



- Z-test 2) Hypothesis Testing

- T-test ✓ H₀, H₁, P value

- CHI square significance value



Ex Let's say there are 20 classes at your college and you have collected the heights of student in the class. Heights are recorded [175cm, 180cm, 160cm, 140cm, 135cm, 160cm, 135cm, 190cm]

Descriptive Question

"What is the average height of the entire classroom?

$$\frac{175 + 180 + 160 + 140 + 135 + 160 + 135 + 190}{8} =$$

Interential Question

"Are the height of the students in classroom similar to what you expect in the entire college?"

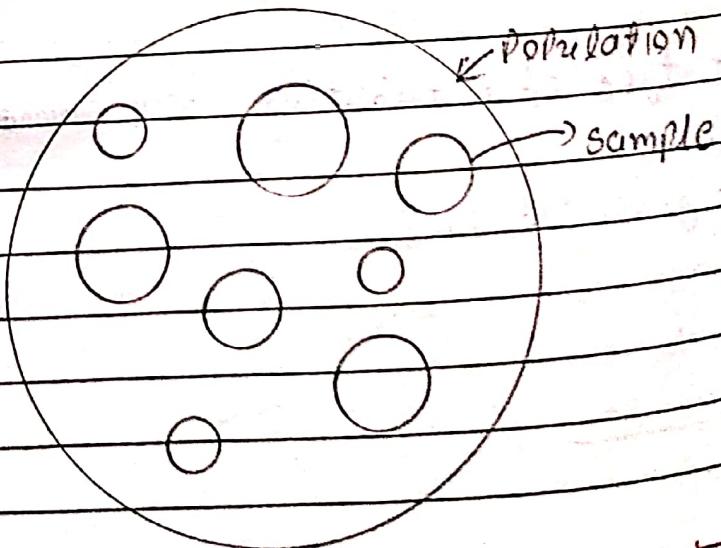
sample

↑ population

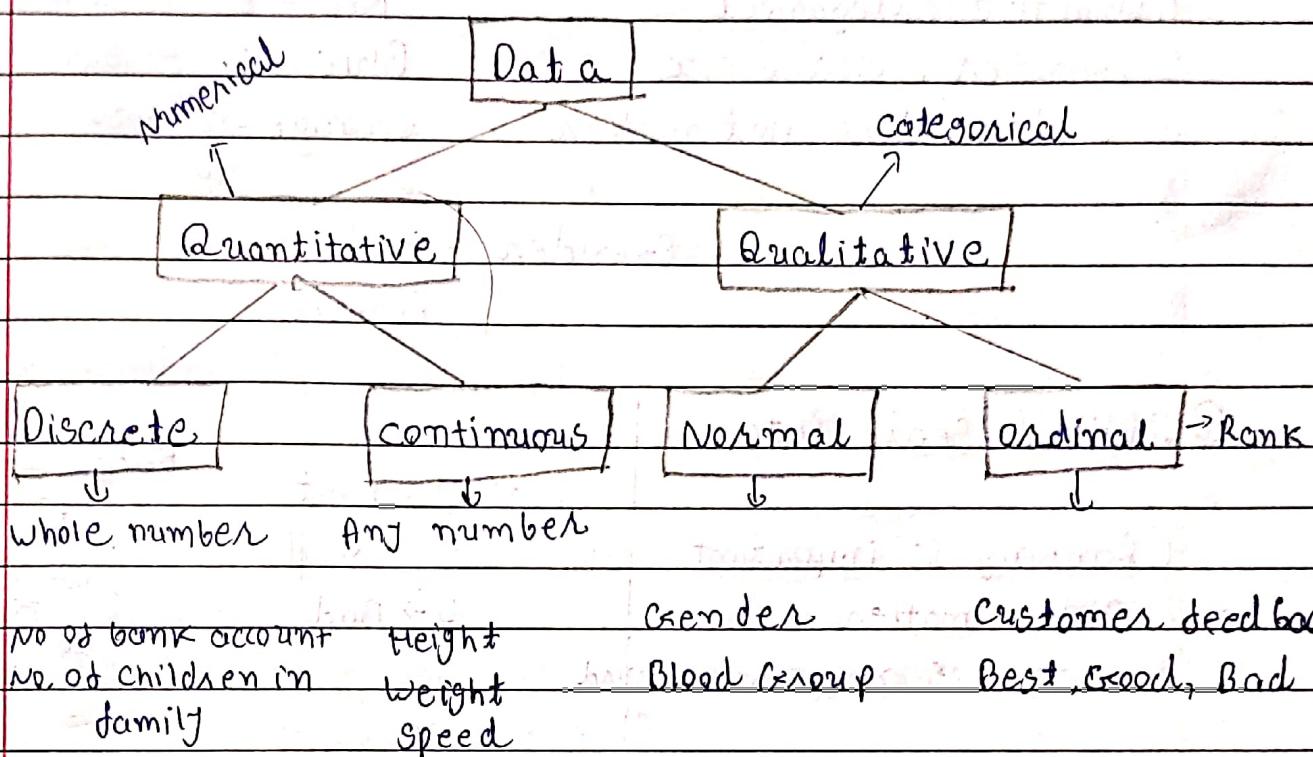
Population and sample data

Population: The group you are interested in studying

Sample: A subset of Population



Types of data



Scale of measurement

1 Nominal Scale Data

2 Ordinal Scale Data

3 Interval Scale Data

4 Ratio Scale Data

① Nominal Scale Data

Eg:- Favourite color

1. Qualitative / Categorical

Red \rightarrow 5 - 501.

2. Gender, color etc.

Blue \rightarrow 3 - 301.

3. Order does not matter

Orange \rightarrow 2 - 201.
 $\frac{10}{10}$

Gender :- Male

Female

② Ordinal Scale Data

Eg:- $\downarrow \rightarrow$ Best

Race

1. Ranking is important

2 \rightarrow Good

1st

2. Order matter

3 \rightarrow Bad

2nd

3. Difference cannot be measured

3rd

③ Interval Scale Data

Eg:- Temperature variable

1. The order matter

30F

OF cm
-30F

2. Difference can be measured

60F

shuttle
pm

3. Ratio cannot be measured

90F

4. No such "0" starting point

120F

④ Ratio Scale Data

Eg:- Student marks in class

1. The order matter

0, 90, 30, 95, 40

2. Difference are measurable

HO

(including ratio)

A.S.S = 0, 30, 40, 75, 90

3. contain a "0" starting point

40 - 30 = 10

50 - 30 = 20

$$\text{Ratio} = \frac{40}{30} = [3:1]$$

Measure of Central Tendency

1. Mean
2. Median
3. Mode

1. Mean :- average of given numbers

Population mean (μ) = $\frac{\sum x_i}{n}$

Eg: $x = \{1, 1, 2, 2, 3, 3, 4\}$

$$\mu = \frac{1+1+2+2+3+3+4}{7} = \frac{16}{7}$$

Sample mean (\bar{x}) = $\frac{\sum x_i}{n}$

$$\bar{x} = \frac{1+1+2+2+3+3+4}{7} = \frac{16}{7}$$

2. Median :- Most frequent value in the sample test

$$x = \{4, 5, 2, 3, 2, 7\}$$

1. Sort the random variable $\{1, 2, 2, 3, 4, 5\}$

2. No. of element $\Rightarrow 6$

3. if count is even

if count is odd

$$\{1, 2, \underline{2, 3}, 4, 5\}$$

$$\{1, 2, 2, \underline{3}, 4, 5, 6\}$$

$$\text{Median} = 3$$

$$\text{Median} = \frac{2+3}{2} = 2.5$$

3. Mode :- Most frequent value in sample test

$$\{2, 1, 1, 1, 4, 4, 5, 6, 6\}$$

$$\text{Mode} = 1$$

Measure of Dispersion [Copy of data]

1. Variance
2. Standard deviation

1. Variance

$$\text{Population variance } \sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$$

$$\text{Sample variance } s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

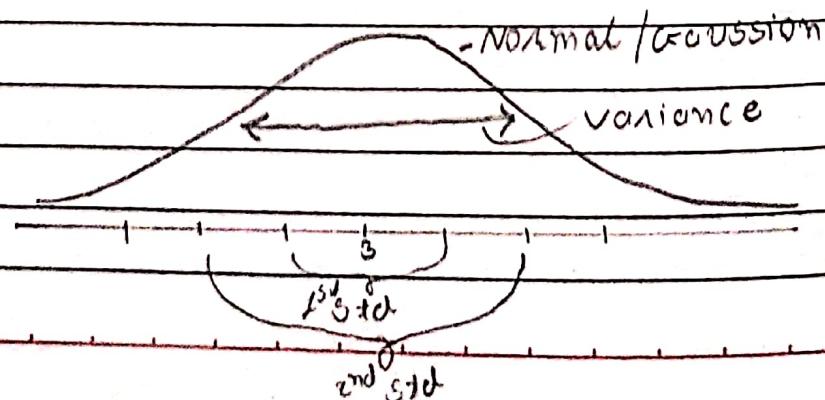
Why we divide Sample Variance by $n-1$?

The sample variance is divided by $n-1$ so that we can create an unbiased estimator of the population variance.

2. Standard deviation

$$\text{Population standard deviation } \sigma = \sqrt{\text{Variance}}$$

$$\text{Sample standard deviation } std = \sqrt{s^2}$$



Random Variable

$$x + 5 = 7 \quad x = 2$$

$$y = j + x \quad j = 6$$

Random variable is a process of mapping the output of a random process or experiments to a number.

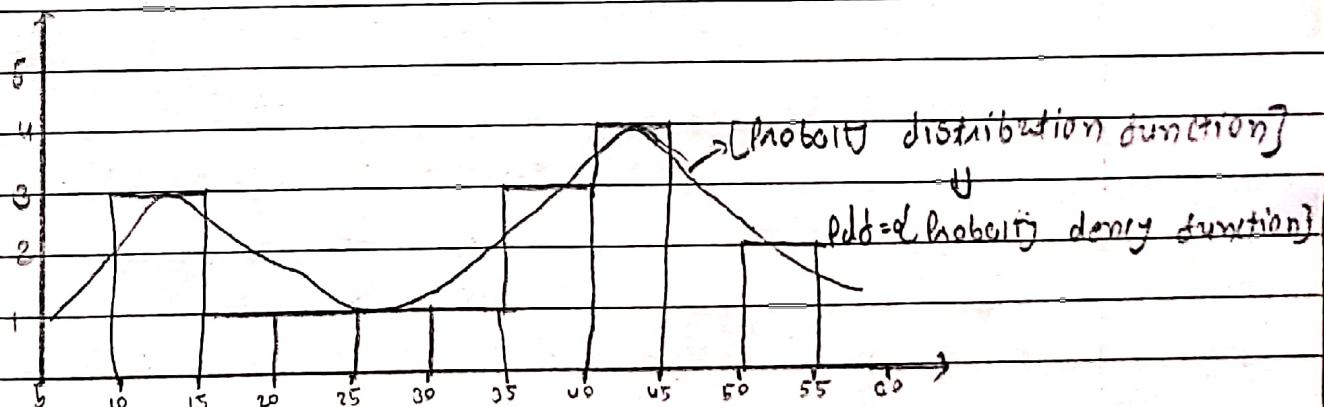
Eg:- Tossing a coin $y = \text{Sum of rolling dice } z \text{ time}$
 $x = \begin{cases} 0 & \text{if Head} \\ 1 & \text{if Tail} \end{cases}$

Histogram and Skewness → [Frequency]

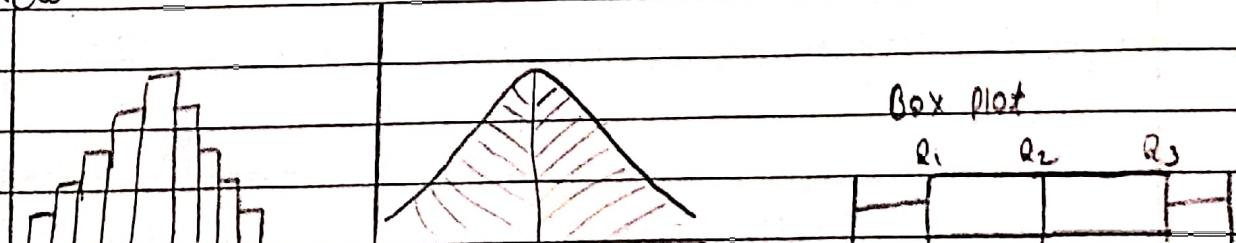
Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51}

$$\frac{50}{10} = 5 \rightarrow \text{bin size } \{ \text{no of bins} = 10 \}$$

$$\frac{50}{20} = 2.5 \rightarrow \text{bin size } \{ \text{no of bins} = 20 \}$$

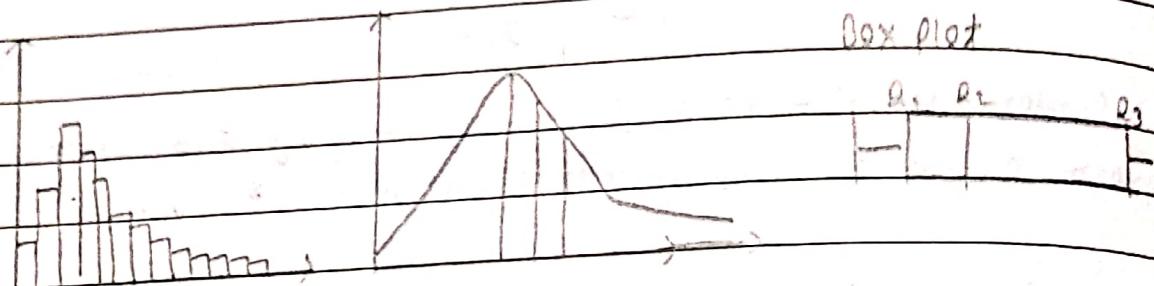


1. Skew

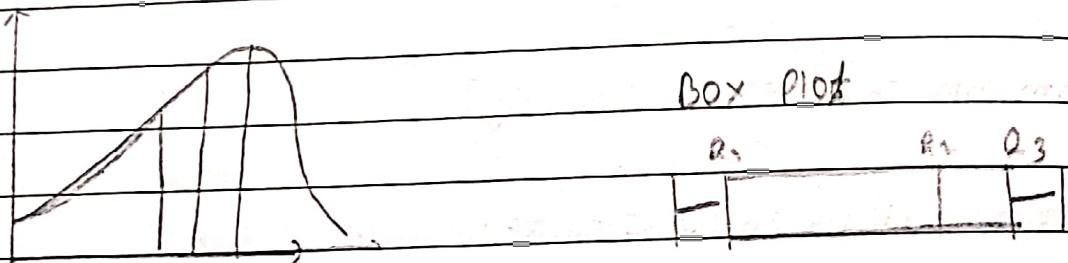


The mean, mode and median all are perfectly at the center

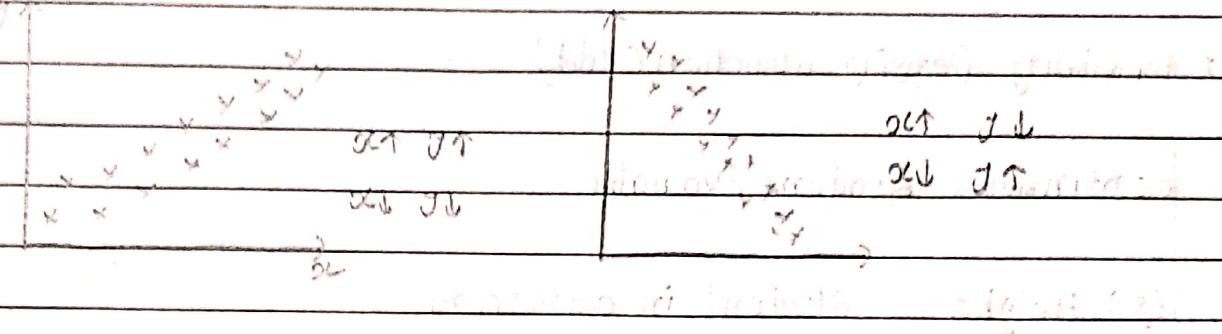
2 Right skewed



3 Left skewed



Covariance & Correlation



Covariance

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

Advantages

Relation between x and y
+ve or -ve value

Disadvantages

covariance does not have a
specific limit value

Pearson correlation coefficient [-1 to 1]

$$g_{x,y} = \frac{\text{cov}(x,y)}{6x \cdot 6y}$$

The more value toward +1 the more +ve correlation it is $g_{x,y}$

The more value toward -1 the more -ve correlation it is $g_{x,y}$

Spearman rank correlation [-1 to 1]

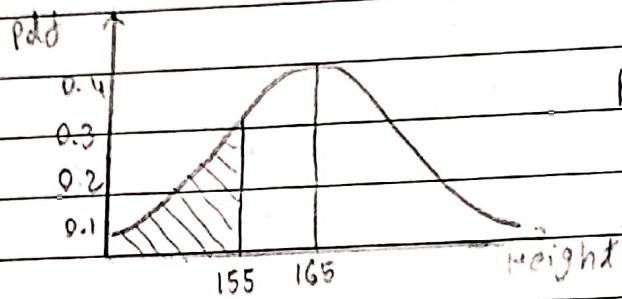
$$g_s = \frac{\text{cov}(R(x), R(y))}{G(R(x)) + G(R(y))}$$

Probability Distribution Function

1. Probability Density function (P.d.f.)

continuous Random Variable

Eg: Height of Student in classroom



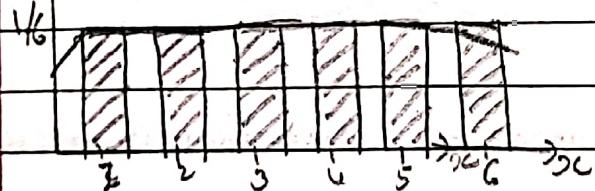
$\Pr(X \leq 155)$ Area under the curve

2. Probability Mass function (P.m.f.)

discrete Random Variable

Eg: Rolling a Dice {1, 2, 3, 4, 5, 6}

Note: (i) If



$$\begin{aligned} \Pr(C_1) &= \frac{1}{6} & \Pr(X \leq 4) &= \Pr(C_1) + \Pr(C_2) + \Pr(C_3) \\ \Pr(C_2) &= \frac{1}{6} & &+ \Pr(C_4) \\ & & &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} \\ & & &= \frac{4}{6} \end{aligned}$$

3. Cumulative Distribution function (C.d.f.)

