

What is Statistics And its application?

→ Statistics is a field that deals with collection, organization, analysis interpretation and presentation of the data.



Decision Making

e.g. Age = {24, 27, 34, 23, 28, 29, 31}

↓ 32 }

{ Online Shopping }

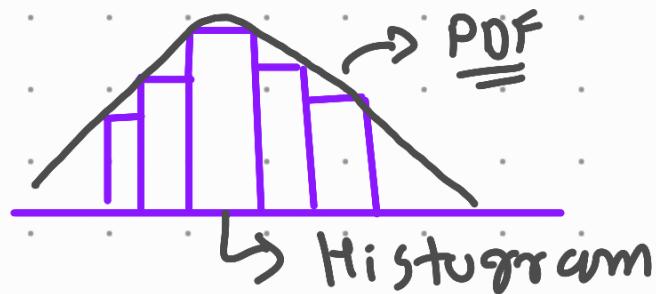
- (1) Mean
- (2) Median of Age
- (3) Distributions of age

→ Create Charts / Graphs

(1) Histogram

(2) PDF

(3) CDF (Cumulative Density Func)

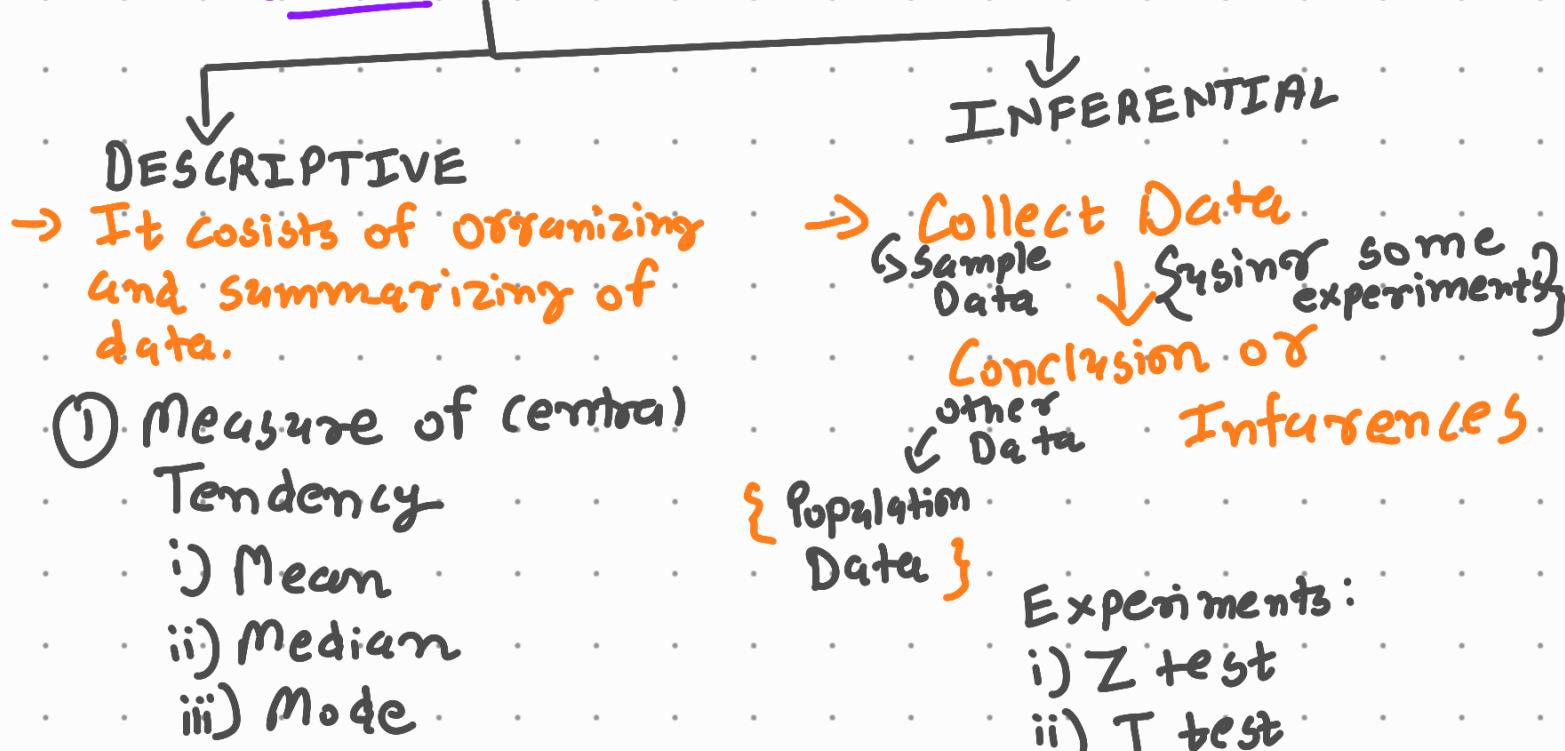


⇒ Final Goal is to understand the data.

Application

- (1) Day to Day usage
- (2) ML / DS
- (3) Data Analyst
- (4) BA
- (5) Risk Analyst
- (6) Research

⇒ Types of Statistics



- (1) Measure of central Tendency
 - i) Mean
 - ii) Median
 - iii) Mode

- (2) Measure of Dispersion
 - i) Variance
 - ii) Standard Deviation

Experiments:

- i) Z test
- ii) T test

e.g. College A \rightarrow 1000 students

class stats

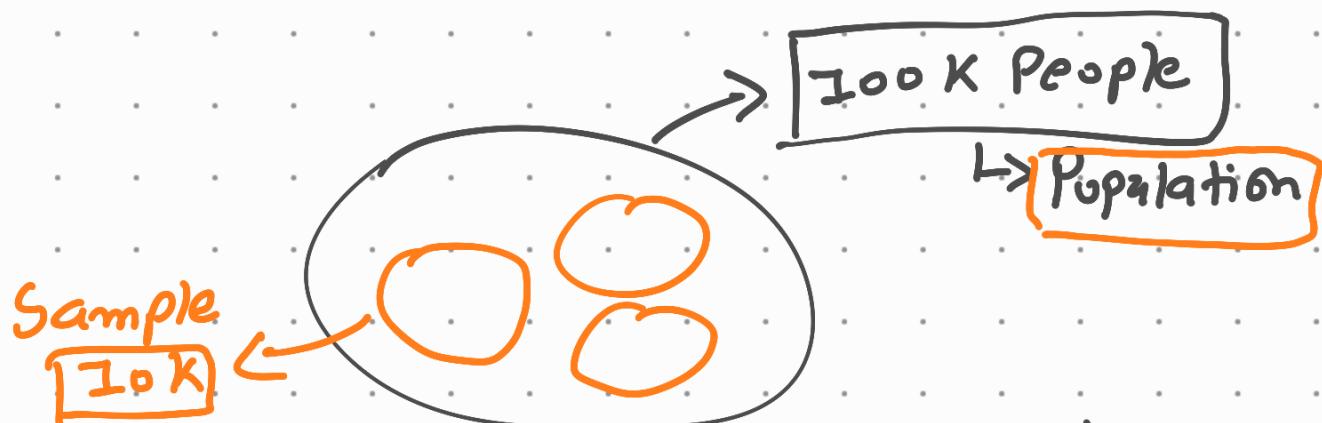
Sample $\rightarrow \{ 180\text{ cm}, 162\text{ cm}, 150\text{ cm}, 160\text{ cm} \}$

Descriptive \rightarrow Mean, Median

Average ≈ 160

Inferential \rightarrow Conclusion, Inferences

\rightarrow Population And Sample



\rightarrow Collect all the weights of all the people
(It is difficult)

\rightarrow Population $\rightarrow N$ } Notation
 \rightarrow Sample $\rightarrow n$ }

e.g. ① Exit Polls

① Measure of Central Tendency

① Mean

② Median

③ Mod

→ Population $\rightarrow N$

→ Sample $\rightarrow n$

① Mean / Avg

$$\mu = \sum_{i=1}^N \frac{x_i}{N}$$

↳ Respect to N

$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n}$$

↳ Respect to n

e.g. Age = {1, 3, 4, 5} \rightarrow Distribution

$$\mu = \frac{1+3+4+5}{4}$$

$$\mu = 3.25$$

② Median

Outlier (Big Number)

$$\text{Ages} = \{1, 3, 4, 5, \underset{\downarrow}{100}\}$$

$$\bar{M} = \frac{1+3+4+5+100}{5} = \frac{113}{5} = 22.6$$

→ To solve the problem of outlier we use Median.

$$\text{Ages} = \{4, 3, 1, 5, 100\}$$

↪ Sort the numbers

$$= \{1, 3, 4, 5, 100\} \rightarrow \text{Median} = 4$$

① Pick the central element
if total number

① Odd → Central element

$$\text{New Ages} = \{1, 3, \boxed{4, 5}, 100, 200\} \rightarrow \text{Median} = 4.5$$

② Even → Both Central elements

$$= \frac{4+5}{2} = 4.5$$

③ Mod

Outlier

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



→ Select the element with the maximum frequency.

O/P $\boxed{= 4} \rightarrow \text{Mod}$

② Measure of Dispersion

① Variance

② Standard Deviation

{ Spread is Less }

$$\text{Ages}_1 = \{ 2, 2, \downarrow 4, 4 \}$$

{ Spread is More }

$$\text{Ages}_2 = \{ 1, 1, \downarrow 5, 5 \}$$

$$\mu = \frac{2+2+4+4}{4}$$

$$\mu = \frac{1+1+5+5}{4}$$

$$\boxed{= 3}$$

Mean ←

$$\boxed{= 3}$$

① Variance

Respect to N,

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

, where N = Population Size

$$AgeS1 = \{2, 2, 4, 4\}$$

$$AgeS2 = \{1, 1, 5, 5\}$$

$$\boxed{\mu = 3}$$

$$\boxed{\mu = 3}$$

$$x_i \quad \mu \quad (x_i - \mu)^2$$

$$2 \quad 3 \quad 1$$

$$2 \quad 3 \quad 1$$

$$4 \quad 3 \quad 1$$

$$4 \quad 3 \quad 1$$

$$\sum = 16$$

$$x_i \quad \mu \quad (x_i - \mu)^2$$

$$1 \quad 3 \quad 4$$

$$1 \quad 3 \quad 4$$

$$5 \quad 3 \quad 4$$

$$5 \quad 3 \quad 4$$

$$\sum = 16$$

$$\sigma^2 = \frac{4}{4} = 1$$

$$\sigma^2 = \frac{16}{4} = 4$$

↳ Dispersion is less.

↳ Dispersion is more

Respect to n ,

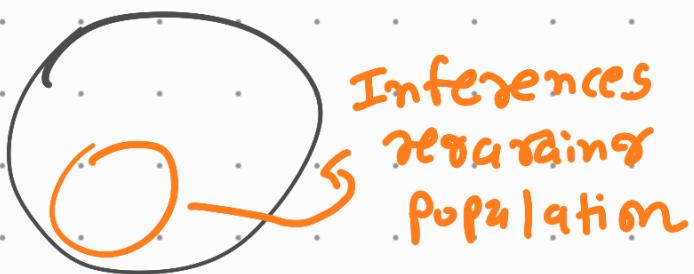
$$S^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$$

, where \bar{x} = Sample mean
 n = Sample Data

→ The main question is why it is divided by $n-1$ instead of n ??

→ This is basically called as Basal Correction.

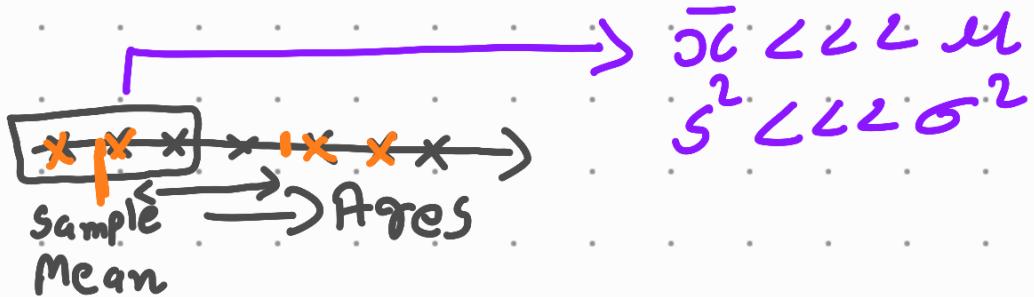
→ Let's consider that currently we use n



Ages = { }

$$\bar{x} \approx \mu$$

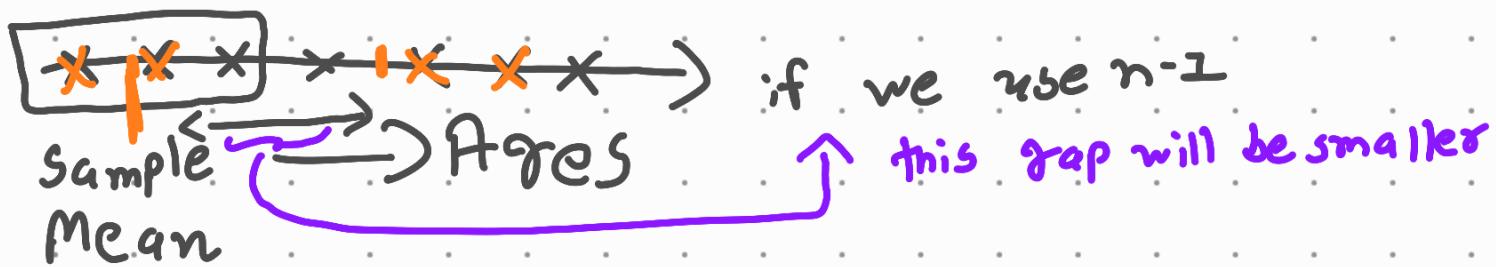
$$S^2 \approx \sigma^2$$



→ let's take

$$s^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}$$

underestimating the true population variance



→ if we use,

$$s^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$$

→ This is known as "Bessel Correction"

→ Then this value is comparatively smaller than n .

→ Also not underestimating the true population variance

$$\rightarrow Dof = n-1$$

Degree of Freedom

② Standard Deviation

→ Respect to N,

$$\text{S.D.} = \sqrt{\sigma^2}$$

→ S.D. Tells that how far a you are from Mean.

→ Respect to n,

$$s = \sqrt{s^2}$$

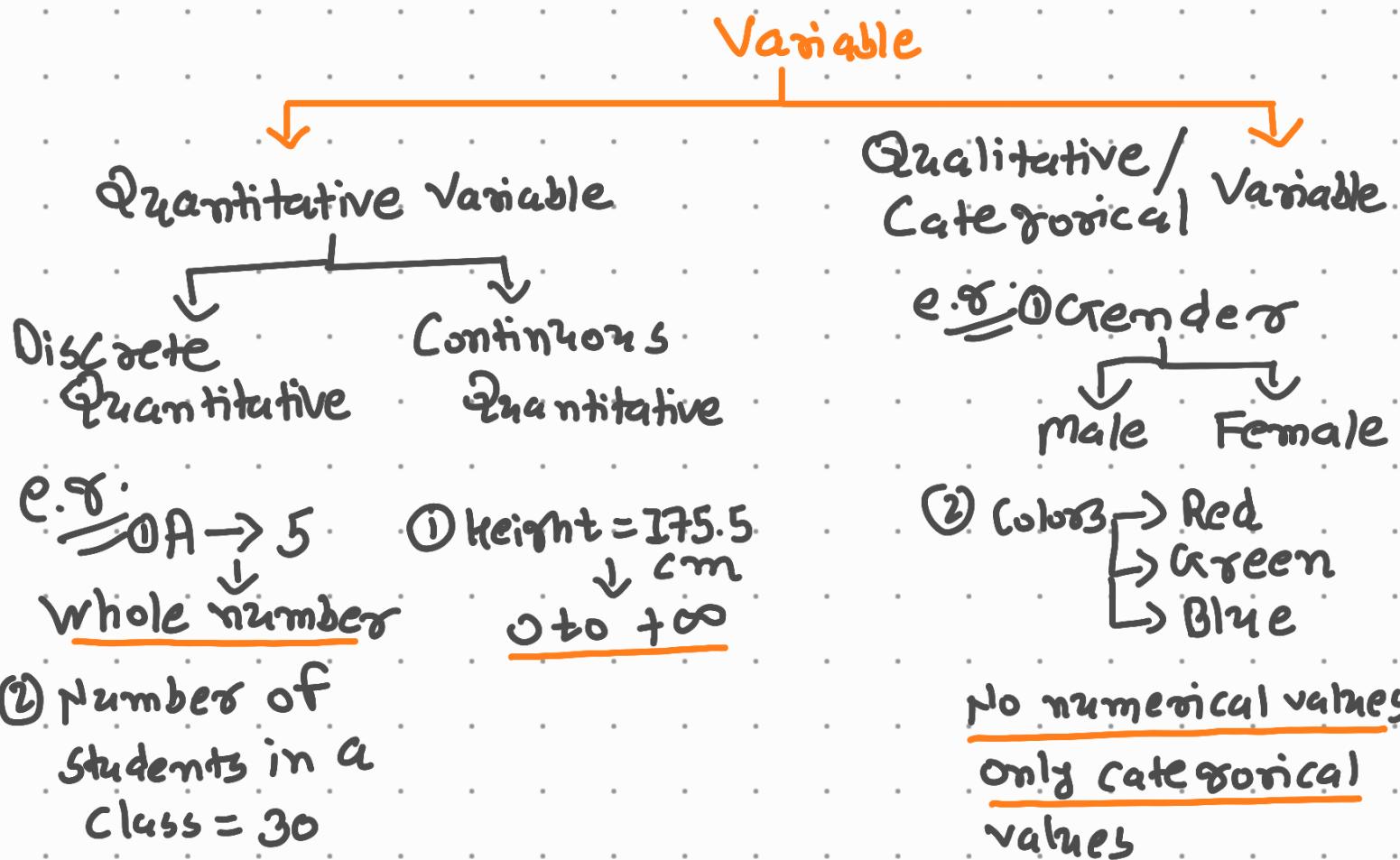
→ What is Variable?

Def:- Variable is a property that can take up any value.

e.g. Age = 25 ✓ Gender = Male ✓

Ages = [12, 20, 25, 36] ✗ Not a variable

→ Different types of Variable



→ Random Variables

Notation - X

→ Random variable is a function whose values are derived from different process or experiments.

$x = 1, x = 2, x = 3$
e.g. $y = 5x + 2$
→ Here x is random variable which can take multiple values.

Random Process

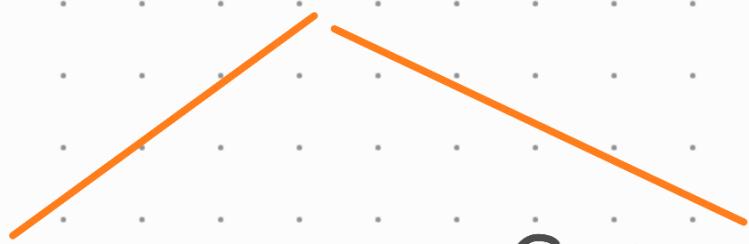
e.g.

$$X = \{ \begin{matrix} 0 & H \\ 1 & T \end{matrix} \text{ Tossing a coin} \}$$

$$X = \{ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} \text{ Rolling a fair dice} \}$$

→ Different Types of Random Variables

Random Variables



Discrete Random variable

e.g. Tossing a coin, }
Rolling a dice }

Continuous Random Variable

e.g. Tomorrow how many inches it is going to rain
↳ 1.1, 5.5, ...

Random Variable

$X \rightarrow$ Function \rightarrow Values

↓
Process or Experiments

Histograms

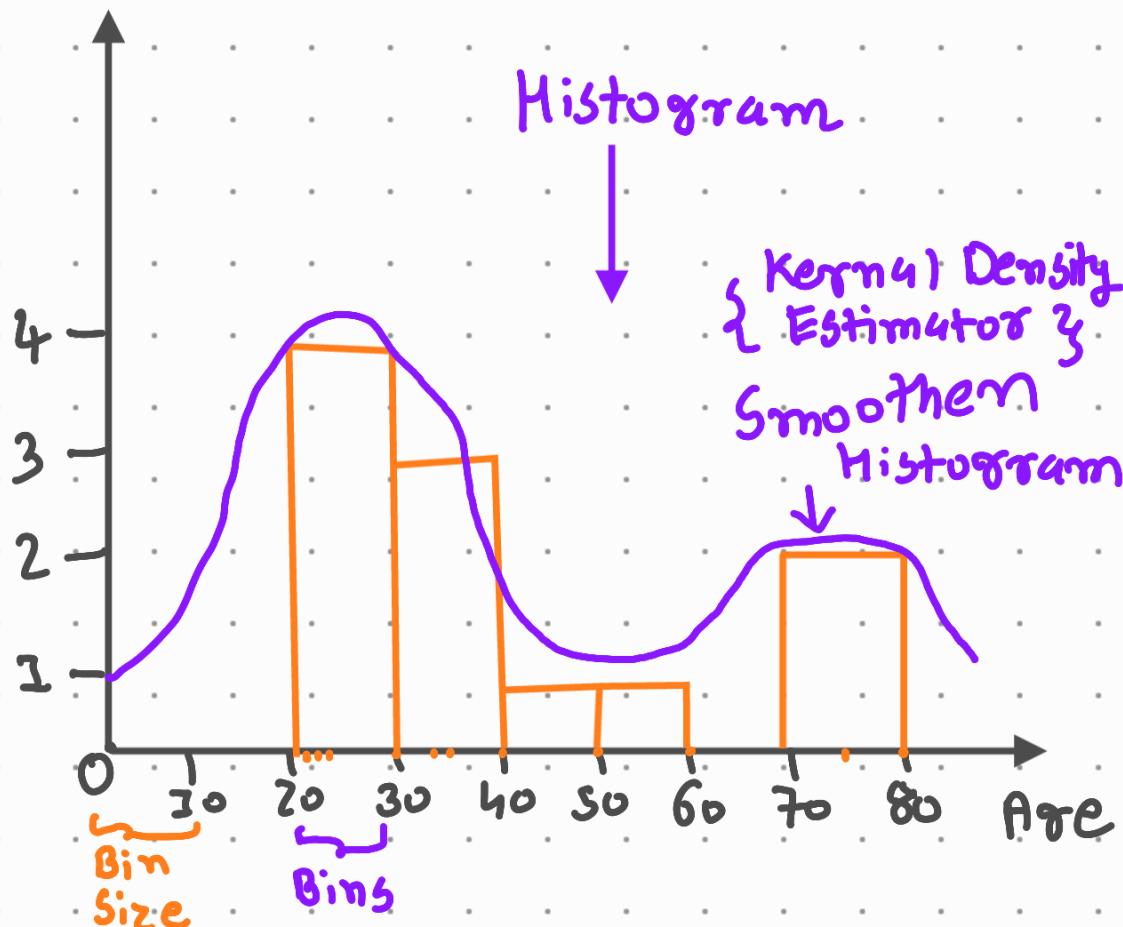
Age

$$X = \{23, 24, 25, 30, 34, 36, 40, 50, 60, 75, 80\}$$

Bins

$$\left\{ \begin{array}{l} 20-30 = 4 \\ 31-40 = 3 \\ 41-50 = 2 \\ 51-60 = 2 \\ 61-70 = 0 \\ 71-80 = 2 \end{array} \right.$$

Count/
Frequency



Percentile And Quartiles

Percentage $\Rightarrow \{1, 2, 3, 4, 5, 6\}$

No of ODD Numbers = 3

Percentage of odd numbers in this group $= \frac{3}{6} \times 100$

$$= 50\%$$

Percentiles: A percentile is a value below which a certain percentage of observations lie.

$\{2, 2, 3, 4, 5, 5, 6, 7, 8, 8, 8, \boxed{9}, \underline{9}, 10\}$

Percentile of value x $= \frac{\text{no of values below } x \times 100}{n}$

$x = 9 \Rightarrow \text{Percentile} = \frac{11}{14} \times 100$

\downarrow
Percentile = 78.57%ile of value 9
Ranking

Value for given $= \frac{\text{Percentile}}{100} \times (n+1)$
Percentile

Value $= \frac{25}{100} \times (15) = 3.75$

\downarrow
Now take $\frac{3+4}{2} = 3.5$

Quartile

25·ile = 1st Quartile

50·ile = 2nd Quartile

75·ile = 3rd Quartile

5 Numbers Summary

- 1) Minimum
- 2) 1st Quartile (25·ile)
- 3) Median
- 4) 3rd Quartile (75·ile)
- 5) Maximum

Outlier

e.g. 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27

[Lower fence \rightarrow Higher fence]

$$\text{Lower fence} = Q_1 - 1.5 \times \text{IQR}$$

$$\text{Higher fence} = Q_3 + 1.5 \times \text{IQR}$$

$$\text{IQR} = Q_3 - Q_1$$

(Inter Quartile Range)

$$Q_1 = \frac{\text{Percentile}}{100} \times (n+1) = \frac{25}{100} (29+1) = 5^{\text{th}} \text{ Position} = 3$$

$$Q_3 = \frac{\text{Percentile}}{100} \times (n+1) = \frac{75}{100} (29+1) = 15^{\text{th}} \text{ Position} = 7$$

$$\text{IQR} = Q_3 - Q_1 = 7 - 3 = 4$$

$$\text{Lower Fence} = 3 - 1.5(4) = -3$$

$$\text{Higher Fence} = 7 + 1.5(4) = 13$$

$[-3 \longleftrightarrow 13]$ → Anything out of this range is considered as an Outlier.

Outlier

1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 12
↳ Removed

→ Minimum = 1

→ 1st Quartile = 3

→ Median = 5

→ 3rd Quartile = 7

→ Maximum = 9

}

5-Number Summary

→ Box Plot (Use to see the Outlier)

Covariance And Correlation

- These are two statistical measures used to determine the relationship between two variables.
- Both are used to understand how changes in one variable are associated with changes in another variable.

① Covariance

Covariance is a measure of the relationship between two random variables. The metric evaluates how much – to what extent – the variables change together. In other words, it is essentially a measure of the variance between two variables. However, the metric does not assess the dependency between variables.

- Positive covariance: Indicates that two variables tend to move in the same direction.
- Negative covariance: Reveals that two variables tend to move in inverse directions.

e.g.

	X	Y
→ 2	3	
→ 4	5	
→ 6	7	
→ 8	9	

[Quantify the relationship between X and Y]

$X \uparrow$	$Y \uparrow$
$X \downarrow$	$Y \uparrow$
$X \uparrow$	$Y \downarrow$
$X \downarrow$	$Y \downarrow$

Y



→ Positive (tive) Covariance



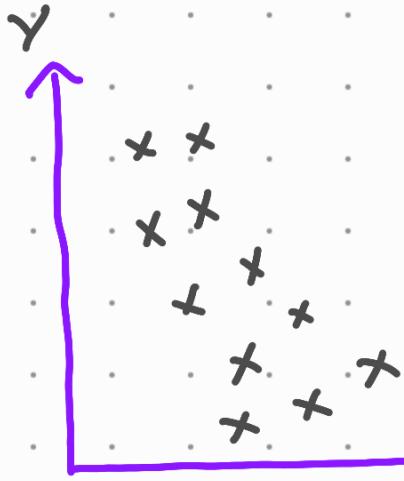
$X \uparrow$	$Y \uparrow$
$X \downarrow$	$Y \downarrow$

→ Dataset

Size of House	Price (Lakhs)
1200	45
1300	50
1500	75

→ Size of House ↑ Price ↑

→ Size of House ↓ Price ↓



$x \downarrow y \uparrow$
 $x \uparrow y \downarrow$

x	y
7	10
6	12
5	14
4	16

\rightarrow Negative (-ve) covariance

\rightarrow Formula

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

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

$$\text{Cov}(x, x) = \text{Var}(x)$$

$\rightarrow x_i \rightarrow$ Data points of Random Variable x

$\rightarrow \bar{x} \rightarrow$ Sample mean of x

$\rightarrow y_i \rightarrow$ Data points of Random Variable y

$\rightarrow \bar{y} \rightarrow$ Sample mean of y

C.8 Students Data

Hour Studied (x)

2
3
4
5
6

Exam Score (y)

50
60
70
80
90

$\rightarrow x \uparrow y \uparrow$ } +ve Covariance
 $\rightarrow x \downarrow y \downarrow$

$$\textcircled{1} \bar{x} = \frac{2+3+4+5+6}{5} = 4$$

$$\textcircled{2} \bar{y} = \frac{50+60+70+80+90}{5} = 70$$

$$\begin{aligned} \text{Cov}(x,y) &= (2-4)(50-70) + (3-4)(60-70) \\ &\quad + (4-4)(70-70) + (5-4)(80-70) \\ &\quad + (6-4)(90-70) \end{aligned}$$

5 - 1

$\text{Cov}(x,y) = 20 \Rightarrow +\text{ve Covariance}$

↓

\rightarrow The number of hours studied increases the Exam score.

e.g.

	x	y
7		10
6		12
5		14

$\rightarrow x \downarrow y \uparrow$ $\rightarrow x \uparrow y \downarrow$ $\{$ -ve Covariance

$$\textcircled{1} \bar{x} = \frac{7+6+5}{3} = 6$$

$$\textcircled{2} \bar{y} = \frac{10+12+14}{3} = 12$$

$$\begin{aligned} \text{Cov}(x, y) &= (7-6)(10-12) + (6-6)(12-12) \\ &\quad + (5-6)(14-12) \\ &= \frac{(1)(-2) + (0)(0) + (-1)(2)}{3-1} \end{aligned}$$

$$= \frac{(-2) + (0) + (-2)}{2}$$

$$\boxed{\text{Cov}(x, y) = 0}$$

Advantages

① Quantify the Relationship between X and Y

Disadvantages

① Covariance doesn't have a specific limit value

$$\text{Cov}(x, y) \Rightarrow -\infty \text{ to } +\infty$$

② Correlation


 Pearson Correlation Coefficient
 Spearman Rank Correlation

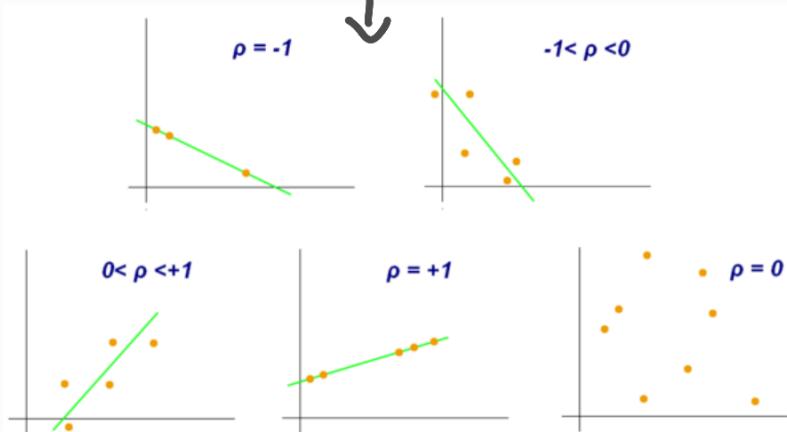
① Pearson Correlation Coefficient

→ It limits the value $\Rightarrow [-1 \text{ to } 1]$

All points in straight line

→ Formula

$$r_{x,y} = \frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y}$$



→ The more the value towards ± 1 , the more tve correlated x & y

→ The more the value towards -1, the more
-ve correlated $x \& y$

C.S. Students Data

Hour Studied (x)

2

3

4

5

6

Exam Score (y)

50

60

70

80

90

→ $x \uparrow y \uparrow \}$ +ve Covariance
→ $x \downarrow y \downarrow \}$

$$\textcircled{1} \bar{x} = \frac{2+3+4+5+6}{5} = 4$$

$$\textcircled{2} \bar{y} = \frac{50+60+70+80+90}{5} = 70$$

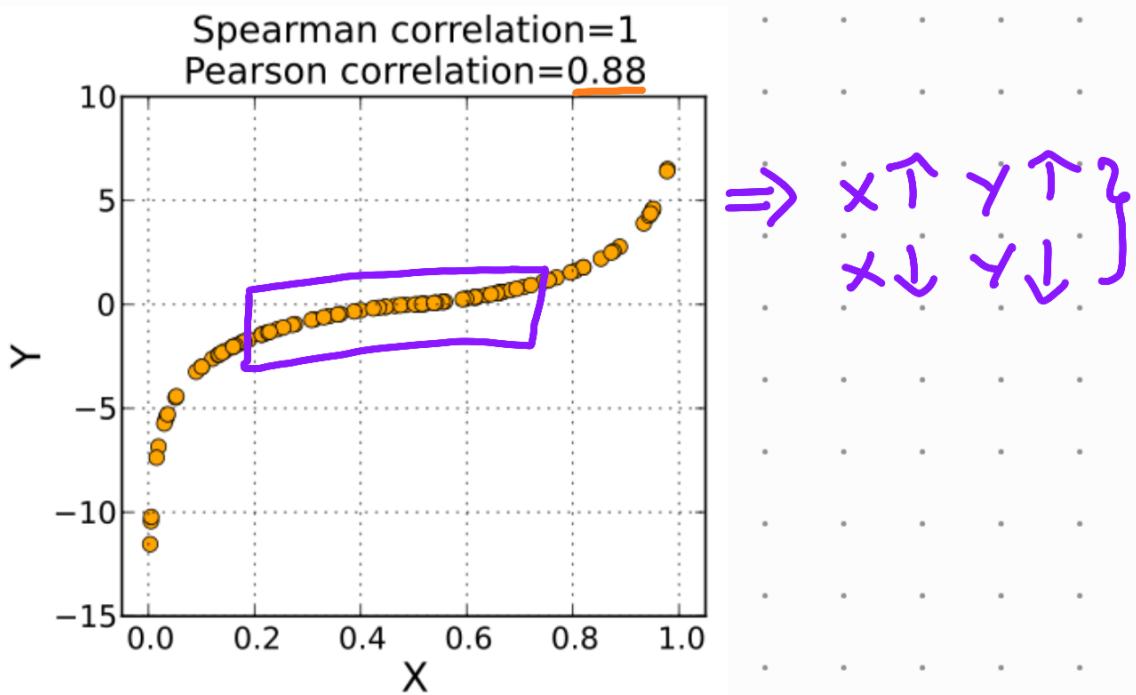
$$\begin{aligned} \text{Cov}(x,y) &= (2-4)(50-70) + (3-4)(60-70) \\ &\quad + (4-4)(70-70) + (5-4)(80-70) \\ &\quad + (6-4)(90-70) \end{aligned}$$

5 - 1

$$\boxed{\text{Cov}(x,y) = 20}$$

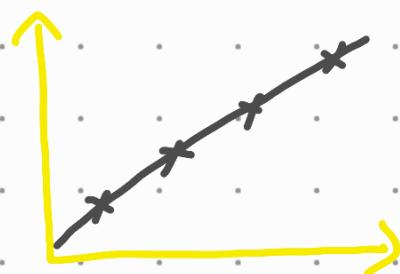
$$S_{x,y} = \frac{20}{\sigma_x \cdot \sigma_y} \Rightarrow 0 \text{ to } 1$$

② Spearman Rank Correlation



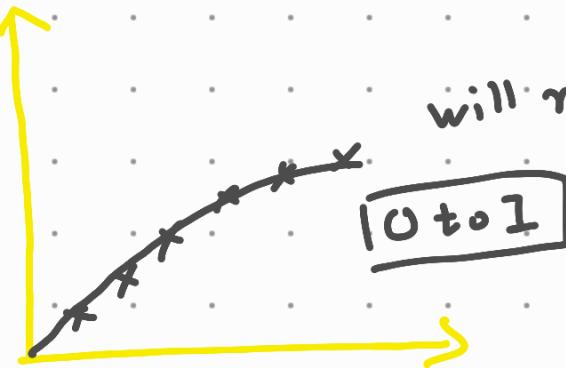
→ Pearson will give -1

if



For Linear Data

if



→ Pearson correlation can't capture for non-linear data.

→ Formula

$$\gamma_S = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) * \sigma(R(y))}$$

R = Rank

e.g.

x	y	R(x)	R(y)
1	2	2	1
3	4	3	2
5	6	4	3
7	8	5	5
0	7	1	4

No effect

Feature Selection

Size of ↑ No. of Room ↑ Location ↑

~~No. of people
staying in
house ↑~~

~~50~~

Haunted

Op →
Price ↑↑↓
-ve correlation

Probability

→ Probability - It is about determining the likelihood of an event.

e.g. ① Toss a coin → {H, T}

$$P_s(H) = \frac{1}{2} = 50\%$$

$$P_s(T) = \frac{1}{2} = 50\%$$

② Rolling a dice → {1, 2, 3, 4, 5, 6}

$$P_s(x=1) = \frac{1}{6}$$

→ Formula

$$P(A) = \frac{\text{Number of Favourable Outcome}}{\text{Total Number of Favourable Outcomes}}$$

→ Mutual Exclusion Events

→ Two events are Mutual Exclusive if they can not occur at the same time.

e.g. ① Tossing a coin



$$P_S(H \text{ or } T) = P_S(H) + P_S(T)$$

$\hookrightarrow \{ \text{Additive Rule for Mutual Exclusive Event} \}$

$$P_S(H) = 1/2 \quad P_S(T) = 1/2$$

$$P_S(H \text{ or } T) = 1$$

(2) Rolling a dice $\rightarrow \{1, 2, 3, 4, 5, 6\}$

$$\begin{aligned} P_S(1 \text{ or } 5) &= P_S(1) + P_S(5) \\ &= 1/6 + 1/6 \end{aligned}$$

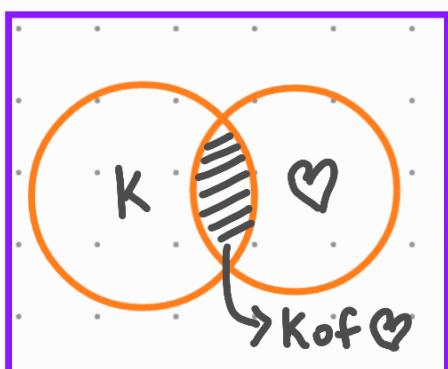
$$P_S(1 \text{ or } 5) = 1/3$$

→ Non-mutual Exclusive Events

e.g. Taking a Card from the deck \rightarrow Total - 52 cards

$$\left(\begin{array}{l} P_S(K \text{ or } \heartsuit) \\ \hookrightarrow \text{Here a card can be a K or } \heartsuit \end{array} \right)$$

$$\begin{aligned} &= P_S(K) + P_S(\heartsuit) - P(K \text{ and } \heartsuit) \\ &= \frac{4}{52} + \frac{13}{52} - \frac{1}{52} \quad \hookrightarrow \{ \text{Additive Rule for Non-mutual Exclusive events} \} \\ &= \frac{16}{52} = \frac{4}{13} \end{aligned}$$



\rightarrow Multiplication Rule $\rightarrow \{$ Independent And Dependent Events $\}$

① Independent Events

\rightarrow 2 events are independent if they do not effect one another.

e.g. ① Tossing a coin $\rightarrow \{ H \text{ and } T \}$

$$P_3(H) = 1/2 \quad P_3(T) = 1/2$$

\rightarrow Getting Tail first time does not effect getting Head then.

② Rolling a dice $\rightarrow \{ 1, 2, 3, 4, 5, 6 \}$

$$P_3(1) = 1/6 \quad P_3(2) = 1/6$$

\rightarrow Both events does not effect each other same as tossing a coin.

② Dependent Events

\rightarrow 2 events are dependent if and only if they effect each other.

e.g. ① Take a King card from the deck and then Queen card from the deck. (without putting cards back)

$$P_S(K) = \frac{4}{52}$$

$$P_S(Q) = \frac{4}{51}$$

↳ These type of events are Dependent Events.

Multiplication Rule Formula

① Independent events

$$P_S(A \text{ and } B) = P_S(A) * P_S(B)$$

② Dependent events

$$P_S(A \text{ and } B) = P_S(A) * P_S(B/A)$$

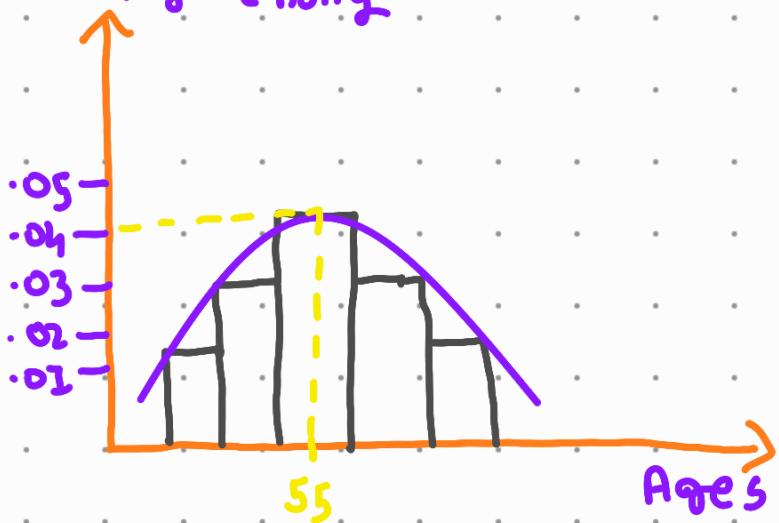
↳ Conditional Probability

Probability Distribution Functions

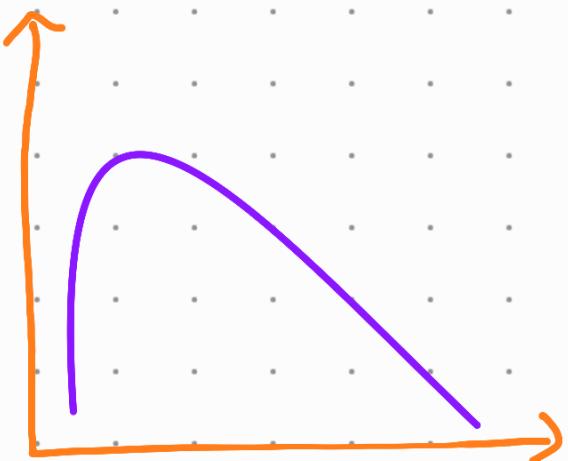
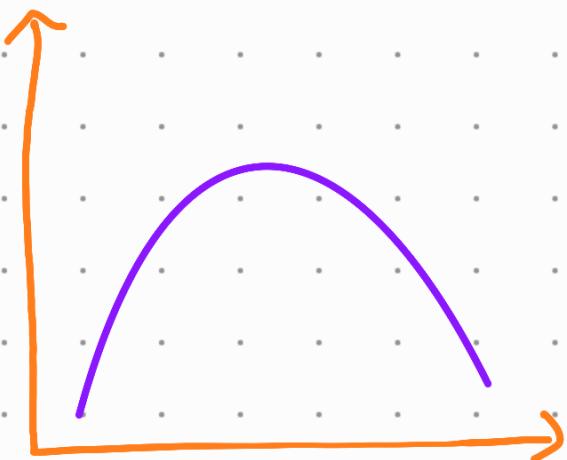
→ PDF describe how the probabilities are distributed over the values of a random variable.

e.g. Ages = {.....} \Rightarrow Continuous Random Variables

Probability Density



→ Variious PDFs



2 Main types of Probability Distribution Function

(1) Probability Mass Functions (PMF)

↳ Used for discrete random variables

(2) Probability Density Function (PDF)

↳ Used for continuous random variables