

Descriptive : Mean, Median,  $\boxed{\text{Avg} = 165} \checkmark$

Inferrential : Conclusion, Inferences }

## Random Variables

$$X \xrightarrow{\text{function}} \text{Values} \quad \begin{cases} y = 5x + 2 & x=1, x=2, n=3 \\ y = 5 \cdot 2 = 7 \end{cases}$$

- Processes or Experiments

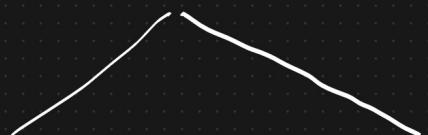
$$X = \begin{cases} 0 & H \vee \text{Tossing a coin} \\ 1 & T \vee \end{cases}$$

$$X = \begin{cases} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{cases} \quad \begin{array}{l} \text{Rolling a } 6\text{-sided} \\ \text{dice} \end{array}$$

Eg: Height of the people attending  
the event tomorrow

Eg: 150cm, 160cm, 160.1cm, ...

## Random Variables



### Discrete Random

Eg: Tossing a coin  
Rolling a dice

### Continuous Random

#### Variable

Eg: Tomorrow how many  
inches it is going  
to rain

[0, 1.1, 5.5, 10.5, 10.75]



# Sample Variance

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

$\uparrow\uparrow\uparrow$

{ Imp interview }      { Population Mean }

{ question }

Sample Mean /  
Population Mean /  
 $\leftarrow$  Inference

$\left\{ \text{Unbiased Estimation} \right\}$

$\left\{ n-2, n-3, n-4 \right\}$

Experimentation     $\left\{ n-1 \right\} \leftarrow \checkmark$

$\underline{\underline{}}$

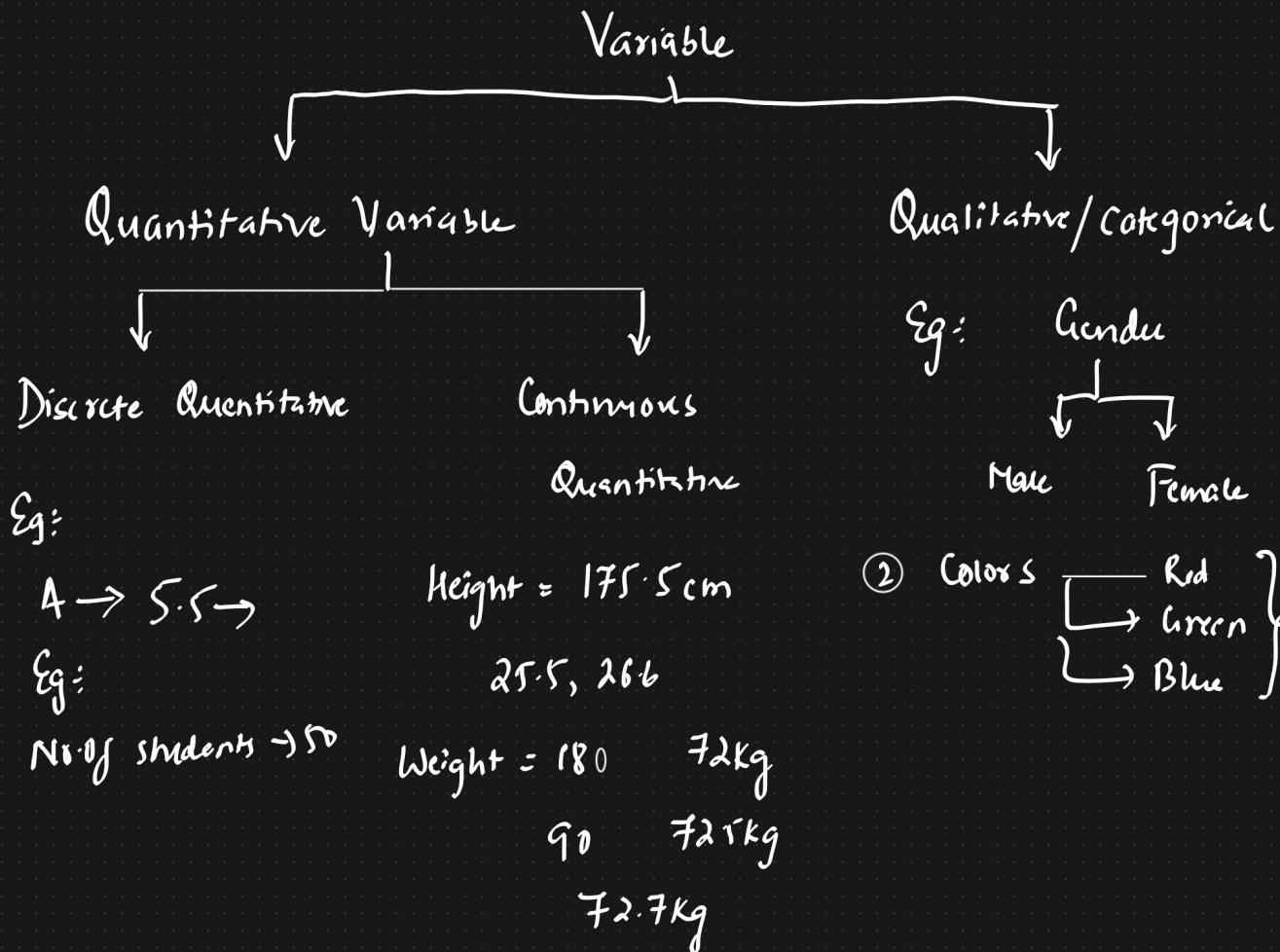
# Variable

Dcfn : Variable is a property that can take up any value

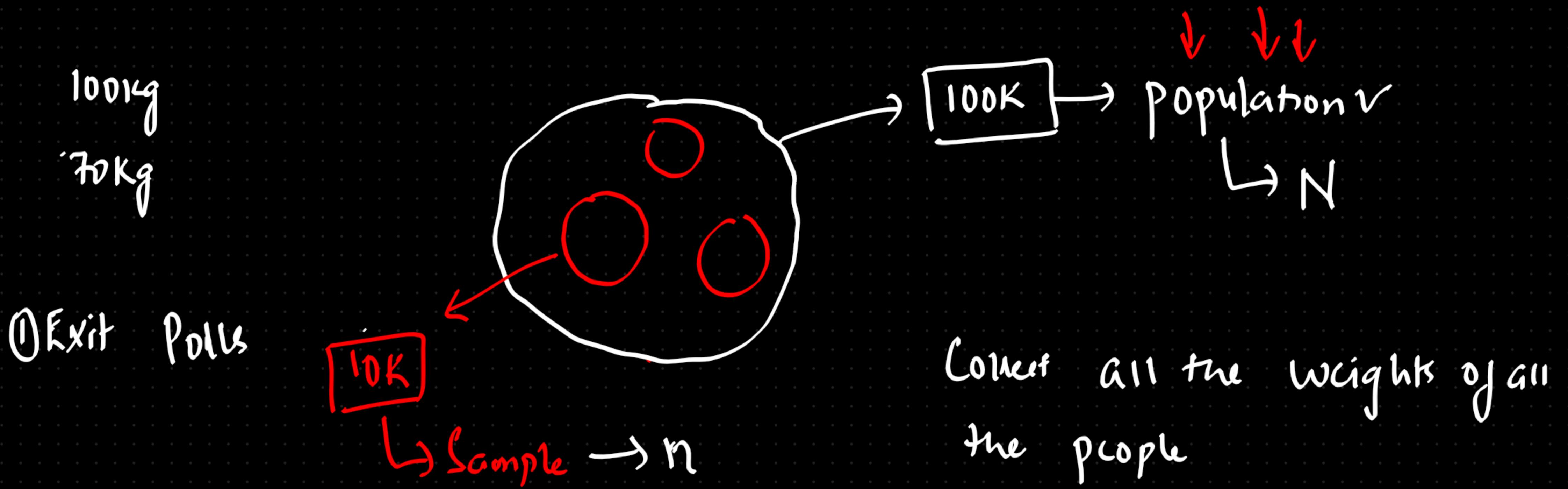
$$\{ \text{Age} = 25 \} \quad \text{Gender} = \text{Male} \quad \text{Height} = 7.2$$

$$\text{Age} = [12, 20, 25, 36, 70] \times \text{Variable}$$

Different types of Variable.



# Population And Sample



# Measure of Central Tendency

① Mean

② Median

③ Mode

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

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

Ages = {1, 3, 4, 5} → Distribution

$$M = \frac{1+3+4+5}{4} = \frac{13}{4} = 3.25$$

② Median

Ages = {1, 3, 4, 5, 100} {Outlier}

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

Ages = {4, 3, 1, 5, 100} → Odd

↳ Sort the numbers {1, 3, 4, 5, 100, 200} → Even

↳ Median → 4

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

Good

# Measure of Dispersion

① Variance

② Standard Deviation

$$Age_1 = \{2, 2, 4, 4\} \rightarrow \text{DISP 1}$$

$$\mu = \frac{2+2+4+4}{4} = 3 \quad \left\{ \begin{array}{l} \text{Spread is less} \\ \text{few} \end{array} \right\}$$

{ Spread is more }

$$Age_2 = \{1, 1, 5, 5\} \rightarrow \text{DISP 2}$$

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

① Variance

Population Data  $\{N\}$  Size

Sample Data  $\{n\}$  → Sample Mean ( $\bar{x}$ )

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

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

$$Age_1 = \{2, 2, 4, 4\}$$

$$Age_2 = \{1, 1, 5, 5\}$$

↓  
Imp Interview

$$\mu = \frac{2+2+4+4}{4} = 3_{II}$$

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

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

$$2 \quad 3 \quad 1$$

$$2 \quad 3 \quad 1$$

$$4 \quad 3 \quad 1$$

$$4 \quad 3 \quad 1$$

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

$$1 \quad 3 \quad 4$$

$$1 \quad 3 \quad 4$$

$$5 \quad 3 \quad 4$$

$$5 \quad 3 \quad 4$$

$$\frac{1}{N} \sum (x_i - \mu)^2 \quad N=4$$

$$N=4$$

$$\sum (x_i - \mu)^2 = 16/4 = 4$$

# Revising All we have learnt.

## Population

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

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

$$\sqrt{\sigma^2} \rightarrow \text{Population Standard deviation}$$

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

## Sample

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

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

$$s = \sqrt{\text{Sample Variance}}$$

↓  
Sample Standard Deviation

# Covariance And Correlation

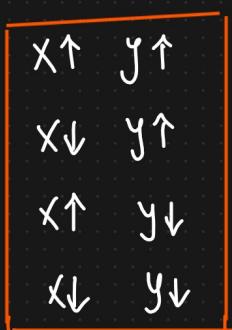
Covariance and correlation are two statistical measures used to determine the relationship between two variables. Both are used to understand how changes in one variable are associated with changes in another variable.

## Covariance

**Definition:** Covariance is a measure of how much two random variables change together. If the variables tend to increase and decrease together, the covariance is positive. If one tends to increase when the other decreases, the covariance is negative.

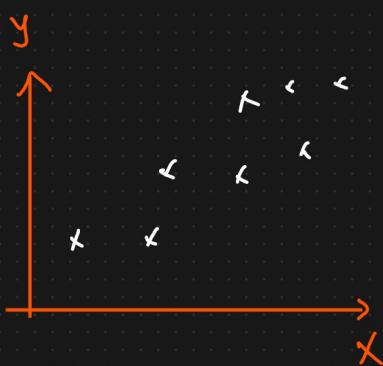
[Quantify the Relationship between X and Y]

X	Y
2	3
4	5
6	7
8	9

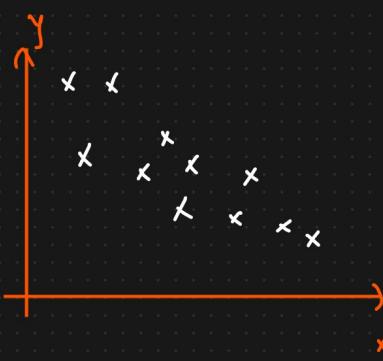


Dataset

↓ ↑ Size of house	Price ↑ ↓
1200	45 lakhs
1300	50 lakh
1500	75 lakh



⇒ +ve Covariance ⇒ +ve value



X	Y
7	10
6	12
5	14
4	16

⇒ -ve Covariance ⇒ -ve value

## Covariance

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

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

$$\boxed{\text{Cov}(X, X) = \text{Var}(X)} \quad \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$$

$x_i \rightarrow$  Data point of random variable  $X$

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

$y_i \rightarrow$  Data points of random variable  $Y$

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

## Students

Hour Studied ( $X$ )

2

3

4

5

6

Exam Score ( $Y$ )

50

60

70

80

90

$x \uparrow y \uparrow \Rightarrow +ve$   
 $x \downarrow y \downarrow$  covariance

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

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

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

$$\text{Cov}(X, Y) = (2-4)(50-70) + (3-4)(60-70) + (4-4)(70-70) + (5-4)(80-70) + (6-4)(90-70)$$

4

$$\text{Cov}(X, Y) = \underline{\underline{20}}.$$

$\Rightarrow$  The positive covariance indicates the no. of hours studied increased the exam score also.

$$\left\{ \begin{array}{l} X \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \quad \begin{array}{l} Y \\ 10 \\ 12 \\ 14 \end{array} \right\} \Rightarrow \underline{\underline{-ve}}$$

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

0.96

0.88

$\text{Cov}(A, B)$

$\text{Cov}(B, C)$

$$\begin{array}{ll} -200 & -300 \\ +100 & +300 \\ \hline 20 & 30 \\ \hline \end{array}$$

-200

+100

20

-300

+300

30

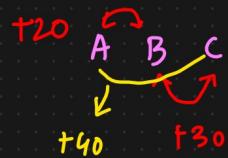
$\text{Cov}(A, B)$

$\text{Cov}(B, C)$

Advantages

[-1 to 1]

Disadvantage



- ① Quantify the Relationship between X and Y

- ① Covariance does not have a Specific limit value.

$$\text{Cov}(X, Y) \Leftarrow -\infty \text{ to } \infty$$

- ② Correlation

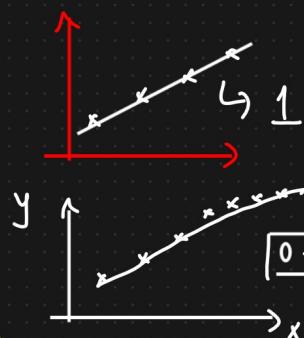
→ Pearson Correlation Coefficient  
→ Spearman Rank Correlation

- ① Pearson Correlation Coefficient  $\Rightarrow [-1 \text{ to } 1]$

$$\rho_{x,y} = \frac{\text{Cov}(X, Y)}{\sigma_x \cdot \sigma_y} = \frac{20}{\sigma_x \cdot \sigma_y} \Rightarrow 0 \text{ to } 1$$

- ① The more the value towards +1 the more +ve correlated X & Y is.  
② The more the value towards -1 the more -ve correlated it is (X, Y)

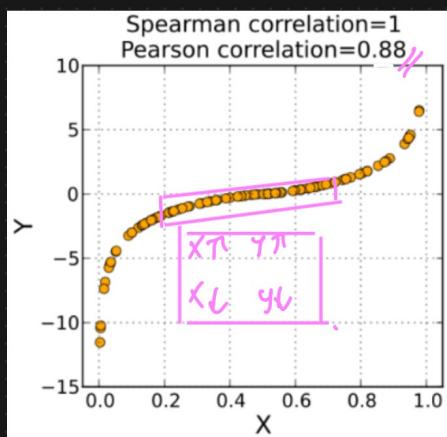
- ③ Spearman Rank Correlation



Pearson Correlation

Correlation for

non linear data



A Spearman correlation of 1 results when the two variables being compared are monotonically related, even if their relationship is not linear. This means that all data points with greater  $x$  values than that of a given data point will have greater  $y$  values as well. In contrast, this does not give a perfect Pearson

$\Downarrow$   
 $\Rightarrow X \uparrow Y \uparrow$   
 $\Rightarrow X \downarrow Y \downarrow$   
 $\Downarrow$

Pearson Correlation  
 $= 0.88$



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

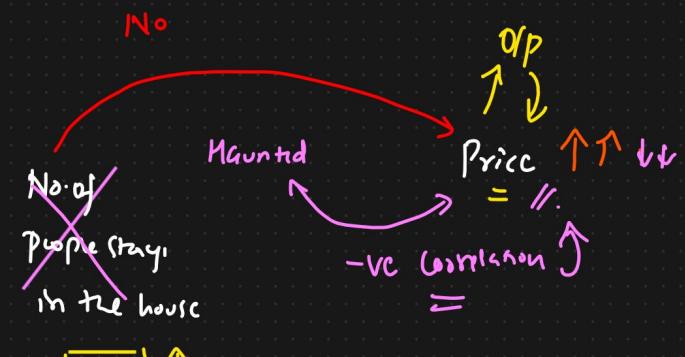
$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) * \sigma(R(y))} \Leftrightarrow$$

## Feature Selection

Size of house  $\uparrow$       No. of Room  $\uparrow$       location  $\uparrow$       ~~No. of people stay in the house~~

are correlated  
=

$\uparrow$



$\boxed{n_o} \uparrow$

# Histograms

of Kernel Density Estimation?

Age

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

20-30 = 4

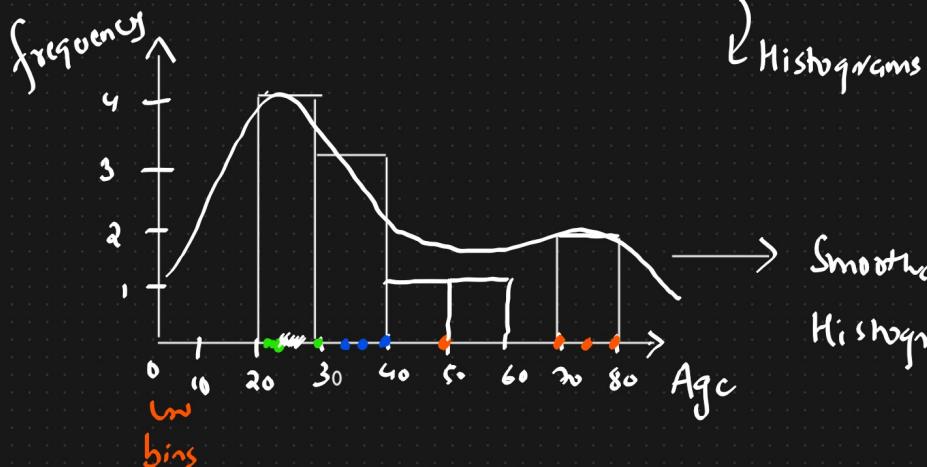
30-40 = 3

41-50 = 1

~~51~~-60 = 1

61-70 = 0

71-80 = 2



Smoothen the  
Histogram



## Percentiles And Quartiles

Percentage = { 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 certain percentage of observations lie.

$$\left\{ \begin{array}{ccccccccc} 2, & 2, & 3, & 4, & 5, & 6, & 7, & 8, & 8, \\ \text{---} & \text{---} \\ 9, & 10 \end{array} \right\} \quad n=9 \quad \frac{3+4}{2} = \underline{\underline{3.5}}$$

$$\text{Percentile of Value } x = \frac{\# \text{ of values below } x \times 100}{n}$$

$$= \frac{11}{14} \times 100$$

$$= 78.57\% \text{ of value } 9$$

Percentile

Ranking

25% is  $\boxed{3.75}$

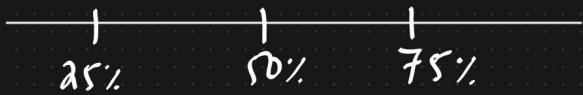
$$\Rightarrow \text{Value} = \frac{\text{Percentile}}{100} \times (n+1)$$

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

$$= \underline{\underline{3.75}} \approx \underline{\underline{3.5}}$$

(2) Quartiles

$25\% = 1^{\text{st}}$  Quartile }  
 $50\% = 2^{\text{nd}}$  Quartile }  
 $75\% = 3^{\text{rd}}$  Quartile }



## 5 Number Summary

1) Minimum ✓ Eg: 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 12 ↗ outlier

2) 1<sup>st</sup> Quartile (25%) [lower fence  $\leftarrow \rightarrow$  higher fence]

3) Median

4) 3<sup>rd</sup> Quartile (75%) lower fence =  $Q_1 - 1.5(IQR)$

5) Maximum (25%) higher fence =  $Q_3 + 1.5(IQR)$

(25%)

$$Q_1 = \frac{\text{Percentile}}{100} \times (n+1) = \frac{25}{100} \times (20) = \underline{\underline{5^{\text{th}} \text{ position}}} = 3 \checkmark$$

(75%)

$$Q_3 = \frac{75}{100} \times (20) = \underline{\underline{15^{\text{th}} \text{ position}}} = 7 \checkmark \quad \text{outlier}$$

$$IQR = Q_3 - Q_1 = 7 - 3 = 4 \checkmark \quad [-3 \longleftrightarrow 13]$$

$$\text{lower fence} = Q_1 - 1.5(IQR)$$

$$= 3 - 1.5(4)$$

$$= 3 - 6 = \boxed{-3}$$

$$\text{higher fence} = Q_3 + 1.5(IQR)$$

$$= 7 + 1.5(4)$$

$$= 13,$$

$$\text{Eg: } 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \boxed{12}$$

Box plot ✓

$$\text{Minimum} = 1$$

$$1^{\text{st}} \text{ Quartile} = 3$$

$$\text{Median} = 5$$

$$3^{\text{rd}} \text{ Quartile} = 7$$

$$\text{Maximum} = 9$$

# Probability

- ① Introduction ✓
- ② Addition Rule (For mutually exclusive event)
- ③ Addition Rule (For non mutually exclusive event)
- ④ Multiplication Rule (Independent & Dependent Events)

{}

① Probability : It is about determining the likelihood of an event

Eg: Toss a coin {H, T}

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

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

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

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

## Mutual Exclusive Event

Two events are Mutual exclusive if they cannot occur at the same time

Eg: Tossing a coin



$$Pr(H) = \frac{1}{2} \quad Pr(T) = \frac{1}{2}$$

$$\begin{aligned} Pr(H \text{ or } T) &= Pr(H) + Pr(T) \quad \{ \text{Additive Rule for mutual Exclusive Event} \} \\ &= \frac{1}{2} + \frac{1}{2} = 1 \end{aligned}$$

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

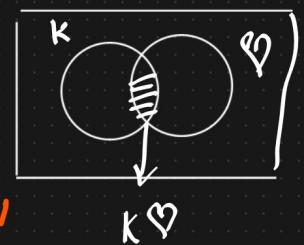
$$\begin{aligned}\Pr(1 \text{ or } 5) &= \Pr(1) + \Pr(5) \\ &= \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}\end{aligned}$$

#### \* Non Mutual Exclusive Events

Eg: Taking a card from the deck

$$\boxed{K} \quad \boxed{K \heartsuit} \quad \boxed{52} \longrightarrow \boxed{K} \text{ or } \boxed{K \heartsuit}$$

$$\begin{aligned}\Pr(K \text{ or } \heartsuit) &= \Pr(K) + \Pr(\heartsuit) - \Pr(K \text{ and } \heartsuit) \\ &= \frac{4}{52} + \frac{13}{52} - \frac{1}{52} \quad \text{Non Mutual} \\ &= \frac{17}{52} - \frac{1}{52} \quad \text{Exclusive Event} \\ &= \frac{16}{52}\end{aligned}$$



# Probability

① Introduction ✓

② Addition Rule (For mutually exclusive event)

③ Addition Rule (For non mutually exclusive event)

④ Multiplication Rule (Independent & Dependent Events)

}

① Probability : It is about determining the likelihood of an event

Eg: Toss a coin {H, T}

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

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

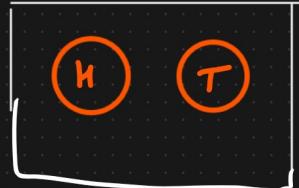
Rolling a dice {1, 2, 3, 4, 5, 6}

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

## Mutual Exclusive Event

Two events are Mutual exclusive if they cannot occur at the same time

Eg: Tossing a coin



$$Pr(H) = \frac{1}{2} \quad Pr(T) = \frac{1}{2}$$

$$\begin{aligned} Pr(H \text{ or } T) &= Pr(H) + Pr(T) \quad \{ \text{Additive Rule for mutual Exclusive Event} \} \\ &= \frac{1}{2} + \frac{1}{2} = 1 \end{aligned}$$

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

$$\begin{aligned}\Pr(1 \text{ or } 5) &= \Pr(1) + \Pr(5) \\ &= \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}\end{aligned}$$

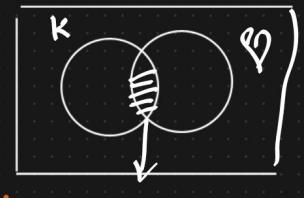
#### \* Non Mutual Exclusive Events

Eg: Taking a card from the deck

$$\boxed{K} \quad \boxed{K \heartsuit} \quad \boxed{5\clubsuit} \longrightarrow \boxed{K} \text{ or } \boxed{K \heartsuit}$$

$$\Pr(K \text{ or } \heartsuit) = \Pr(K) + \Pr(\heartsuit) - \Pr(K \text{ and } \heartsuit)$$

$$= \frac{4}{52} + \frac{13}{52} - \frac{1}{52} \quad \begin{array}{l} \hookdownarrow \\ \text{Non Mutual} \\ \text{Exclusive Event} \end{array}$$



$$= \frac{17}{52} - \frac{1}{52}$$

$$\Pr(K \text{ or } \heartsuit) = \frac{16}{52}$$

#### Multiplication Rule {Independent And Dependent Events}

2 events are Independent if they do not affect one another

Eg: Tossing a coin  $\{H \text{ and then Tails}\}$

$$\Pr(H) = \frac{1}{2} \quad \Pr(T) = \frac{1}{2}$$

Eg: Rolling a dice

$$\Pr(1) = \frac{1}{6} \quad \Pr(2) = \frac{1}{6}$$

## Dependent Events

2 Events are Dependent if they affect each other

Eg: Take a King from the deck and then the Queen Card from the deck

$$\cdot \quad 52 \quad \rightarrow K$$

$$Pr(K) = \frac{4}{52} \quad Pr(Q) = \frac{4}{51}$$

## Multiplication Rule

① Independent Event {Tossing a Coin}

$$Pr(H \text{ and } T) = Pr(H) * Pr(T)$$

$$= \frac{1}{2} * \frac{1}{2}$$

$$\boxed{52} = \frac{1}{4}$$

② Dependent Event  $\nearrow$  Conditional Probability

$$Pr(K \text{ and } Q) = P(K) * Pr(Q/K)$$

$$= \frac{4}{52} * \frac{4}{\cancel{51}}$$

$$= \cancel{\cancel{\cancel{\quad}}}$$