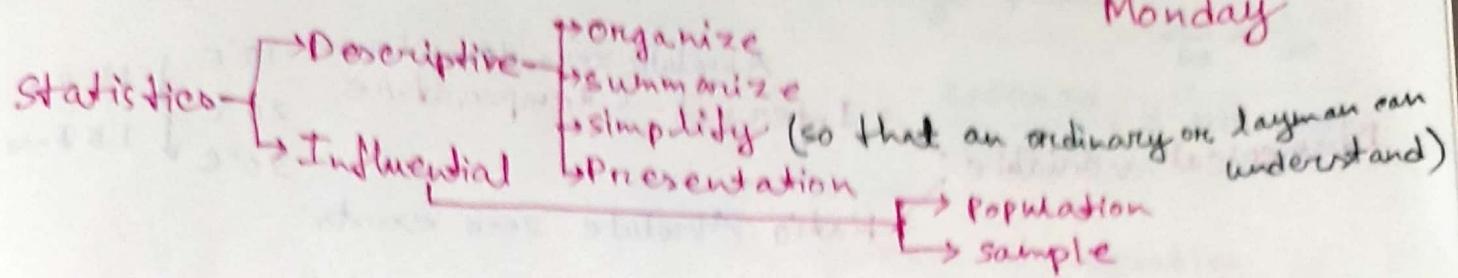


06.03.2023

Monday



Variable / Parameter / Features.

Categorical variable (Location) → Nominal, Ordinal

Numerical variable (Salary)

We convert all categorical variable into numerical variable when using in machine learning.

Nominal: Gender (male/female), location, nationality (No order)

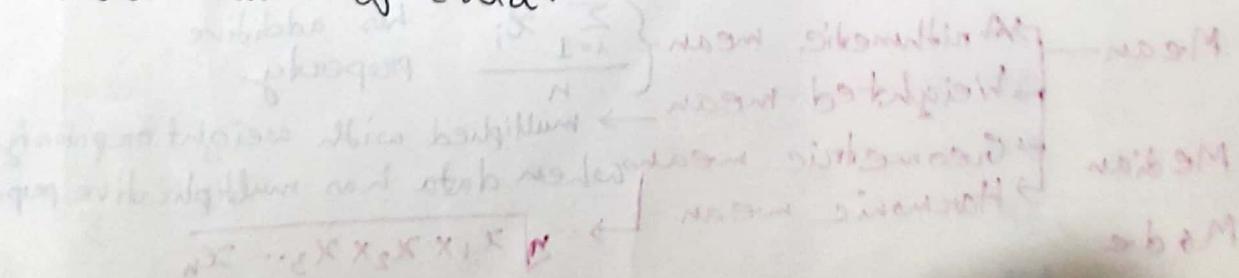
Ordinal: Agreement (strongly disagree, disagree, agree etc.), rating, frequency (always, often, sometimes etc.) (can be ordered)

Univariate: When study is based on a single variable/feature.

→ Multivariate: When study is based on multiple variable/feature.

Often used

Visualization of data.



13.03.2023

Monday

Numerical variable

- Interval → temperature  
Absolute zero doesn't exist
- Ratio Absolute zero exists

$10^{\circ}\text{C}$  ↓ 1.5 times  
 $20^{\circ}\text{C}$  ↓  
 $30^{\circ}\text{C}$  ↓ 1.5 times

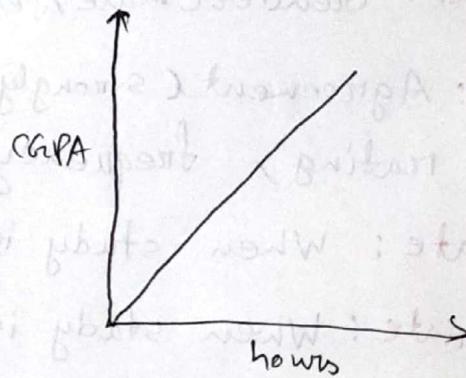
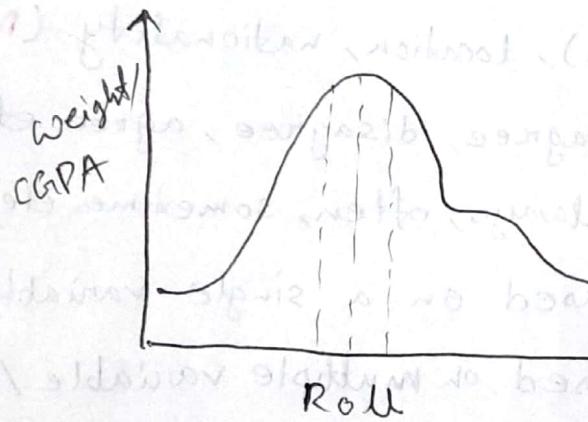
Statistics : Descriptive statistics (works) & Inferential statistics (inference)

Inferential statistics (inference) (works) & (impacts)

Descriptive measures : Central Tendency Measures

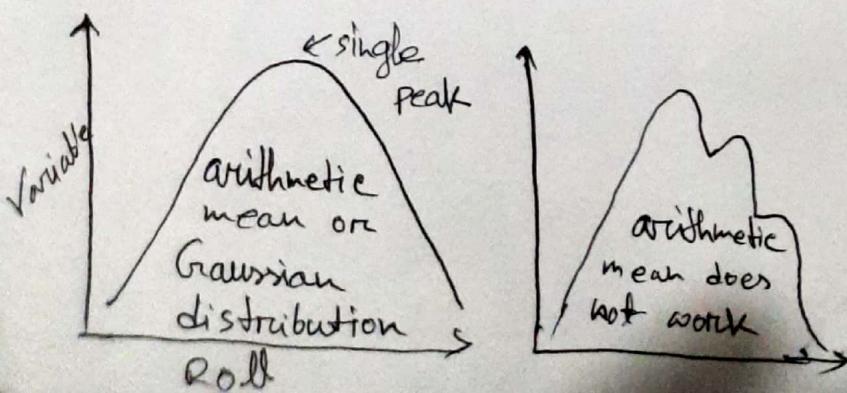
Variation or Variability Measures

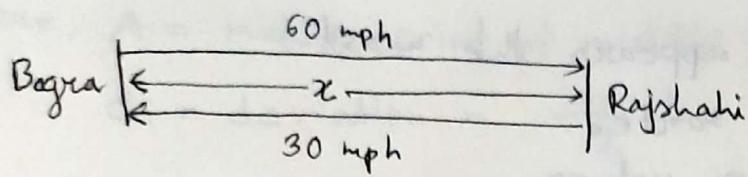
Relative Standing Measures



Measures of Central Tendency (or Location) :

- Mean
  - Arithmetic mean  $\left\{ \frac{\sum_{i=1}^n x_i}{n} \right\}$  when data has additive property
  - Weighted mean → multiplied with weight or priority
- Median
- Mode
- Geometric mean → when data has multiplicative property
- Harmonic mean  $\sqrt[n]{x_1 x_2 x_3 \dots x_n}$





$$\text{Average speed} = \frac{2}{\frac{1}{60} + \frac{1}{30}}$$

$$x = v_1 t_1 = 60 t_1; t_1 = \frac{x}{60} \quad \therefore 2x = v(t_1 + t_2) = \frac{2x}{60 + 30} = 40 \text{ mph}$$

$$x = v_2 t_2 = 30 t_2; t_2 = \frac{x}{30} \quad v = \frac{2x}{\frac{x}{60} + \frac{x}{30}} = 40 \text{ mph}$$

$$\text{Harmonic mean} = \frac{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}{n}; n = \text{number of variables}$$

14.03.2023

Tuesday

Book: Introduction to probability and statistics - Schaum's outline

The median: number of element / odd, median =  $\frac{n+1}{2}$   
 number of element / even, median =  $\frac{(n_2)+(n_2+1)}{2}$

[data needs to be sorted]

It can be computed for

Ratio-level, Interval-level, Ordinal-level,  
 open ended frequency distribution if the median does not lie in an open-ended class.

Height (cm)	Number of girls	
	frequency	cumulative frequency
< 140	4	4
140 - 145	7	11
145 - 150	18	29
150 - 155	11	40
155 - 160	6	46
160 - 165	5	51

$$n = 51$$

$$C_f = 11$$

$$n_2 = 25.6 \quad f = 18$$

$$\text{median} = l + \left( \frac{\frac{n}{2} - C_f}{f} \right) \times h$$

∴ 145 - 150 group → median class

$$h = 5 \quad d = 145$$

**The Mode:** Value that appears the most used when data is not sorted.  
not affected by extreme values  
There may be several modes or no mode (when unique)

Mean - Mode = 3(Mean - Median) used for average less skewed distribution



Skewed distribution      Normal distribution

15.03.2023

Wednesday

**Mean: Grouped data:**

Range	frequency( $f_i$ )	class middlepoint ( $x_i$ )	$f_i x_i$	$d_i$
0 - 10	7	5	=	-2
10 - 20	8	15	=	-1
20 - 30	20	25 A	=	0
30 - 40	10	35	=	1
40 - 50	5	45	=	2

$$\text{Mean, } \bar{x} = \frac{\sum f_i x_i}{\sum f_i} \leftarrow \text{Direct method}$$

$$\text{Mean, } \bar{x} = A + \frac{\sum f_i d_i}{N} + h \leftarrow \text{Assumed mean method}$$

$$= A + \frac{\sum f_i d_i}{\sum f_i} + h$$

Here,  $A$  = middle point of middle class (assumed mean)

$$d_i = \text{deviation} = \frac{x_i - A}{h}$$

$h$  = class size

Mode: Grouped data:

Range	frequency
30 - 40	18
40 - 50	37
50 - 60	45
60 - 70	27
70 - 80	15
80 - 90	8

$$\text{Mode} = L + \frac{\Delta_1}{\Delta_1 + \Delta_2} \times h = 50 + \frac{8}{8+18} \times 10 = 53.08$$

Here,  $L$  = lower limit of mode class

$\Delta_1$  = frequency of mode class - frequency of previous class

$\Delta_2$  = frequency of mode class - frequency of next class

$h$  = class size

Measure of dispersion or measure of spread:

I) Range (maximum data - minimum data)

II) Standard deviation /  $\sqrt{\text{variance}}$

$$\text{Variance} = \frac{\sum (x - \bar{x})^2}{n} \leftarrow \text{for population}$$

$$\text{Variance} = \frac{\sum (x - \bar{x})^2}{n-1} \leftarrow \text{Sample}$$

$S_1$	$0 x_1$
$S_2$	$100 x_1$

$$\bar{x} = 50$$

$$\text{variance} = \frac{(0-50)^2 + (100-50)^2}{(n-1)} = \frac{2500 + 2500}{9} = 5000$$

$$\text{or variance} = \frac{\sum x_i^2 - \bar{x}^2}{n} = \frac{50^2 + 50^2}{4} = 5000$$

Standard deviation,  $SD = \sqrt{\text{Variance}}$

for normal distribution,  $SD = +1 \rightarrow +3$  range

20.03.2023

## Measure of dispersion:

### 1. Range

$$R = L - S$$

Absolute measure of variation  
L = Largest value      Relative measure of variation  
S = Smallest value

Limitation: cannot tell us anything about the character in the distribution within two extreme observations.

### 2. Interquartile Range / Quartile deviation

### 3. Mean deviation

Ungrouped data:  $M.D = \frac{\sum |x_i - \bar{x}|}{N}$

Grouped data:  $M.D = \frac{\sum f_i |x_i - \bar{x}|}{N}; \bar{x} = A + \frac{\sum f_i d_i}{N} \times h$

every data is scattered from mean in (mean deviation)

range.

4. Variance: Ungrouped data:  $\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{N}$

(derivation) Prove  $\rightarrow = \frac{x_i^2 - (\bar{x})^2}{N}$

Grouped data:  $\sigma^2 = \frac{\sum f_i (x_i - \bar{x})^2}{N}$

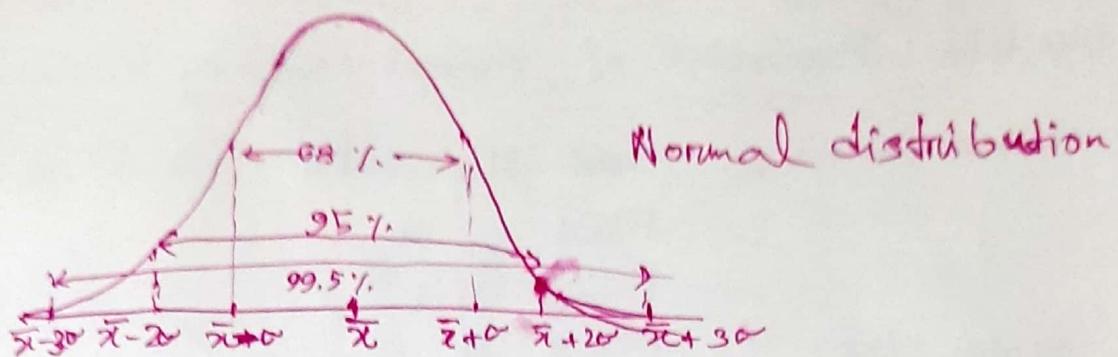
(derivation) Prove  $\rightarrow = \left[ \frac{\sum f_i d_i^2}{N} - \left( \frac{\sum f_i d_i}{N} \right)^2 \right] \times h^2$

5. Standard deviation:  $\sigma = \sqrt{\text{Variance}}$

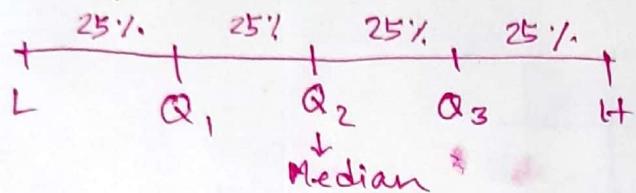
Ungrouped data,  $\sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}}$

Monday

Grouped data:  $s = \sqrt{\frac{\sum f_i (x_i - \bar{x})^2}{n}}$



Interquartile range:



Example:

odd: 2 3 5 7 8 9 10 12 15  
 $\downarrow$   
 $\downarrow$   
 $\downarrow$   
 $Q_1$        $Q_2$        $Q_3$

even: ~~2 3 5 7 8~~ 10 12 14 15 14 16 17 18 10 19 17 17

10. 10  $\underline{12} \quad 14$  14  $\underline{15} \quad 16$  17  $\underline{17} \quad 18$  19  
 $\downarrow$   
 $Q_1 = 13$        $Q_2 = 15.5$        $Q_3 = 17$

$$Q_i = L + \frac{\frac{i \times N}{4} - P.c.f}{f} \times h ; i = 1, 2, 3$$

Deciles:  $D_i = L + \frac{\frac{i \times N}{10} - P.c.d}{f} \times h$

Percentile:  $P_i = L + \frac{\frac{i \times N}{100} - P.c.f}{f} \times h$

22.03.2023

Wednesday

Midquartile: Parameter of central tendency measure.

Five Number Summary: Sort the data

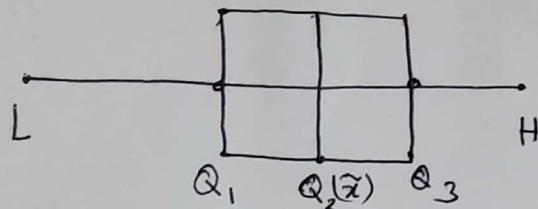
Find L, H

Find  $Q_2(\tilde{x})$

Find  $Q_1, Q_3$

Box whisker plot:

Graphical representation:



Percentile:

$$\left[ \begin{array}{c|c} \text{at most } k\% & \text{at most } (100-k)\% \\ \hline L & P_k & H \end{array} \right]$$

Ungrouped data:  $I = n \frac{P}{100}$   $\rightarrow$   $P^{\text{th}}$  percentile  
position of the percentile  $\downarrow$  Sample size

Sort the data: 11 11 14 15 16 16 17 19 22 25 26 27 31 34 36

$$n = 15$$

$$P = 30$$

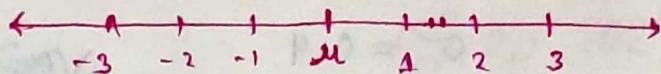
$$I = 15 \frac{30}{100} = 4.5 \approx 5 \quad \rightarrow \text{round it up to next integer}$$

$\therefore$  Percentile at 5<sup>th</sup> position  $\rightarrow 16$

When  $I = n$  (integer) Percentile will be  $\frac{n^{\text{th}} + (n+1)^{\text{th}}}{2}$

Z-score: (standard score)

$$(z = -3 \text{ to } +3)$$



Metric  
Metric

to describe data

Class A

$$\mu = 80$$

$$\sigma = 5$$

Marks = 87

$$z = \frac{87 - 80}{5} = 1.4$$

Class B

$$\mu = 73$$

$$\sigma = 6$$

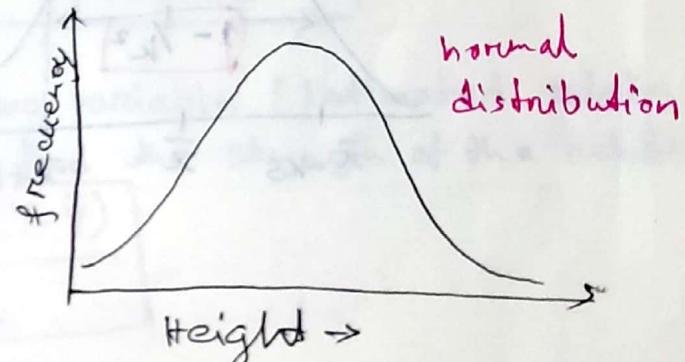
Marks = 82

$$z = \frac{82 - 73}{6} = 1.5$$

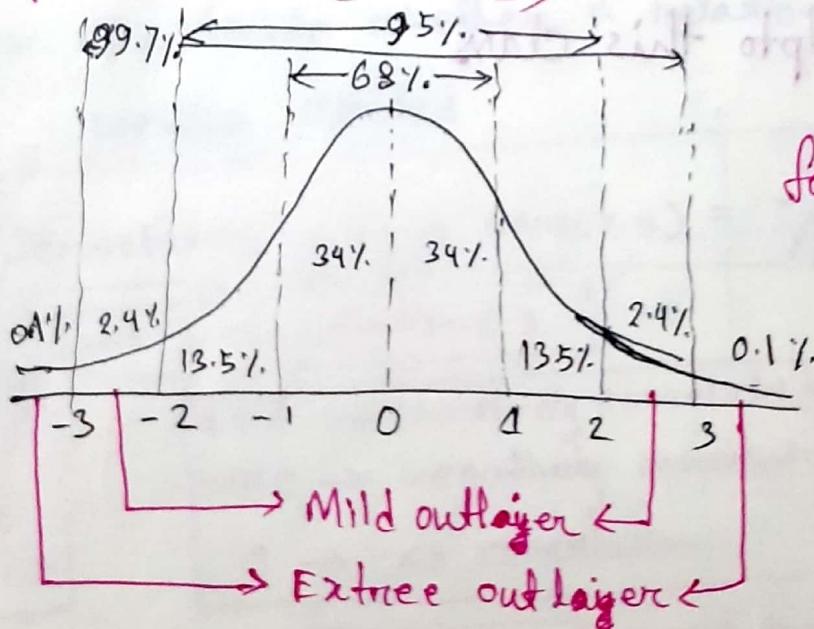
(better performance)

Frequency distribution:

Distribution of male height:



The Empirical Rule (z-score)



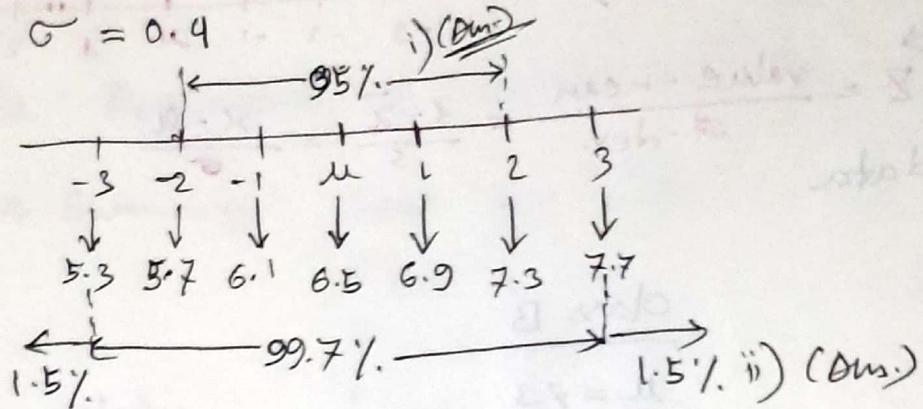
Outlier: data with peculiar characteristics

Standardization  
Standardized data

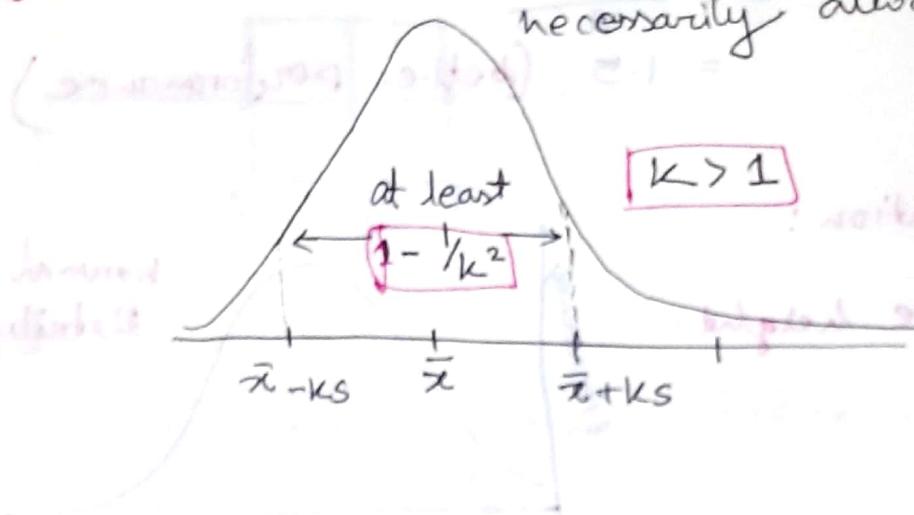
Example:  $\mu = 6.5$

$$\sigma = 0.4$$

i), ii):



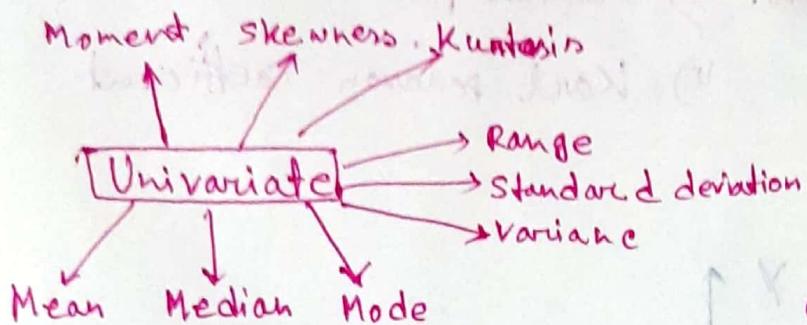
**Chebyshov's Theorem:** When the distribution is not necessarily always normal.



CT syllabus: Upto this class

Wednesday

02.05.2023



Stock price of company X & Y:

DAY	X	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$
1	30	5	-3 ↓	-1 ↓	3
2	35	8	+2 ↑	+2 ↑	4
3	40	8	+7 ↑	+2 ↑	14
4	25	4	-8 ↓	-2 ↓	16
5	35	5	+2 ↑	-1 ↓	-2
Mean	33	6			sum = 35

$$\text{Covariance} = \frac{35}{4}$$

Covariance: Relation between two variable (Estimated relation is not ideal) (cannot find the strength of the relation)

$$\text{Covariance} = \frac{\sum (x - \bar{x})(y - \bar{y})}{(x, y) n - 1}$$

Correlation: Can decide whether a relation is strongly or weakly related.

$r = \text{Karl Pearson Coefficient}$

$$\text{Correlation } (x, y) = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

↑ standard deviation

works only on linear relation

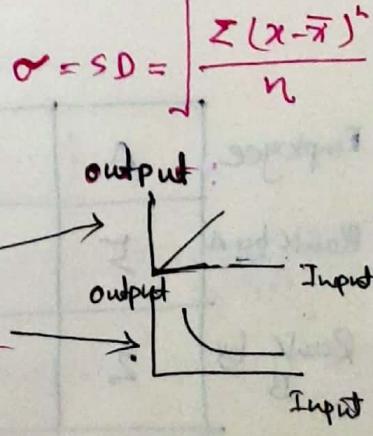
$$-1 \leq \text{correlation} \leq 1$$

+ve → positively correlated  
-ve → negatively correlated  
0 → no correlation

$r = \text{Karl Pearson Coefficient}$

Positive  
Negative

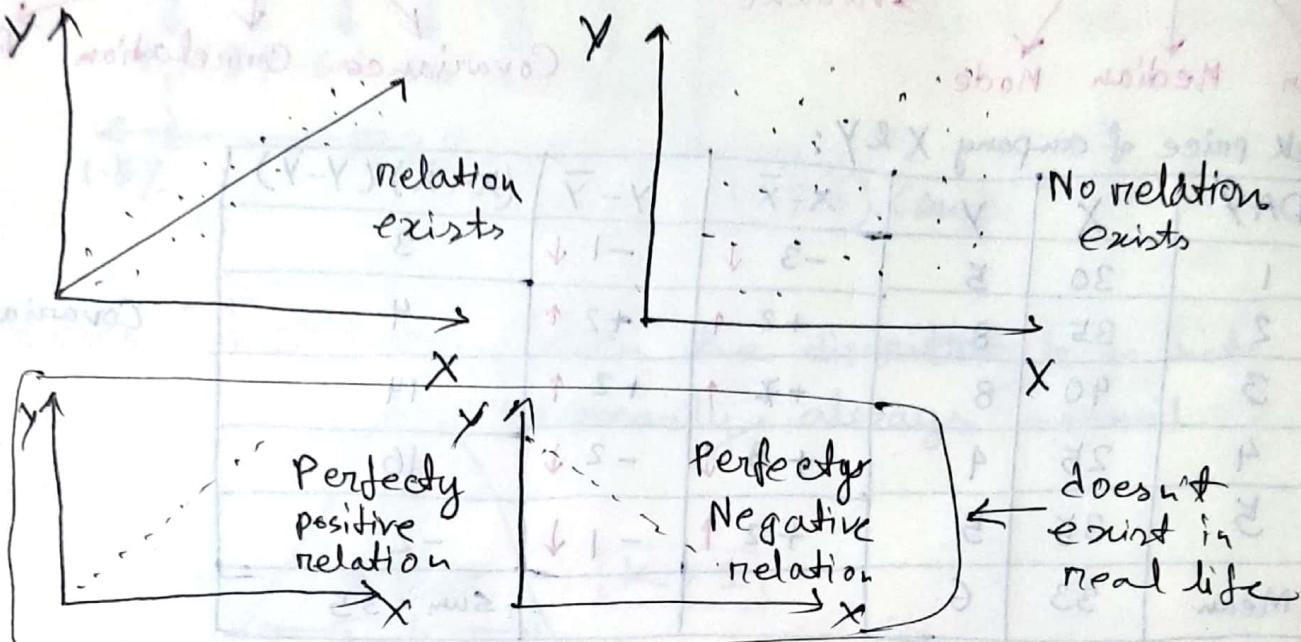
Correlation → Linear  
Non linear  
Simple Multiple  
two variable depend on each other  
one variable depends on multiple variables



# Simple - Linear Correlation:

- Scatter plot
- Karl Pearson coefficient

## Scatter plots:

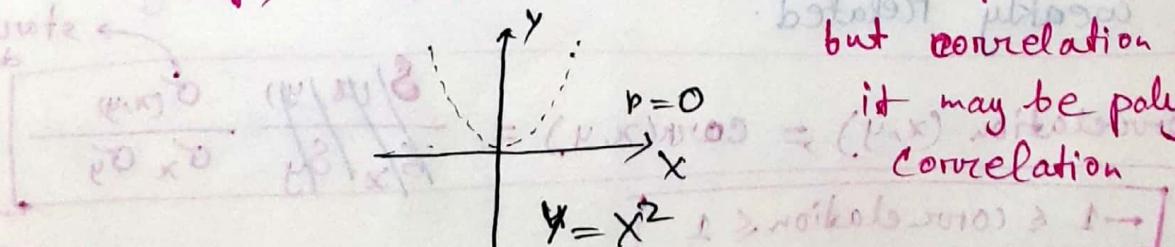


## Karl Pearson coefficient:

$$r = \frac{\sigma_{x,y}}{\sigma_x \sigma_y} = \frac{(x-\bar{x})(y-\bar{y})}{\sqrt{\sum(x-\bar{x})^2} \sqrt{\sum(y-\bar{y})^2}}$$

if,  $r=0 \rightarrow$  there is no simple-linear correlation,

but correlation exists and  
it may be polynomial correlation



Employee	A	B	C	D	E	F
Rank by A	5	4	3	2	1	6
Rank by B	2	4	5	6	7	8

Special grade for  
Karl Pearson  
coefficient for  
ordinal data:

Spearman's  
Rank Algorithm

## Spearman's Rank Algorithm formula!

For unique rank values  $\rightarrow$

$$P_s = 1 - \frac{6 \sum D^2}{N(N^2 - 1)}$$

$$D^2 = (R_1 - R_2)^2$$

Monday

08.05.2023

Example: Ranking relation between Manager 1 & Manager 2

Employee	$R_1$	$R_2$	$D^2 = (R_1 - R_2)^2$
A	10		
B	2		
C	1		
D	4		
E	3		
F	6		
G	5		
H	8		
I	7		
J	9		
$N = 10$			

$$r_s = ?$$

Example: Compute rank correlation coefficient

To rank: sort the preliminary test and rank in ascending order.

Tie in rank (in the same column):

$$r_s = 1 - \frac{6 \left\{ \sum D^2 + \frac{1}{12} (m_1^3 - m_1) + \frac{1}{12} (m_2^3 - m_2) + \dots \right\}}{(N^3 - N)}$$

Example:

Marks in Bangla

	$R_1$	Fraction	$R_2$	$D^2 = (R_1^2 - R_2^2)$
15	2	40	6	
20	3.5	30	9	
28	5	50	1	
12	1	30	9	
40	6	20	2	
60	7	10	1	
20	3.5	30	4	
80	8	60	8	
$N=8$				

Marks in English

Fraction

$$D^2 = (R_1^2 - R_2^2)$$

$$\begin{array}{ccccccccc}
 12 & 15 & 20 & 20 & 28 & 40 & 60 & 80 & 10 & 20 & 30 & 30 & 30 & 30 & 40 & 50 & 60 \\
 R_1 \rightarrow 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & R_2 \rightarrow 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\
 R_1 \rightarrow 1 & 2 & \sqrt{3.5} & \sqrt{3.5} & 5 & 6 & 7 & 8 & R_2 \rightarrow 1 & 2 & \sqrt{3} & \sqrt{4} & \sqrt{5} & 6 & 7 & 8 \\
 m_1 \rightarrow (2) & & 3+4 & 3+4 & & & & & m_2 \rightarrow (3) & & 3+4+5 & 3+4+5 & 3+4+5 & & & \\
 \end{array}$$

$r_s = 0 \rightarrow$  No correlation between Marks in Bangla & Marks in English

Statistics

Descriptive  
(Correlation)

relationship  
exists  
mean  
descriptive

Inferential  
(regression)  
(classification)

Predict  
Continuous  
Value

Regression Analysis:

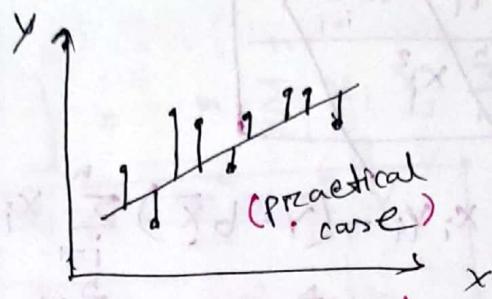
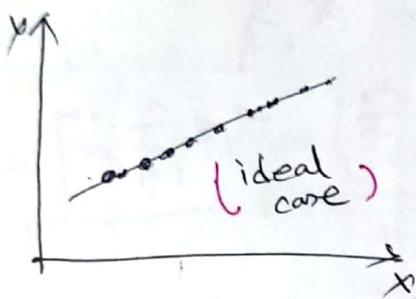
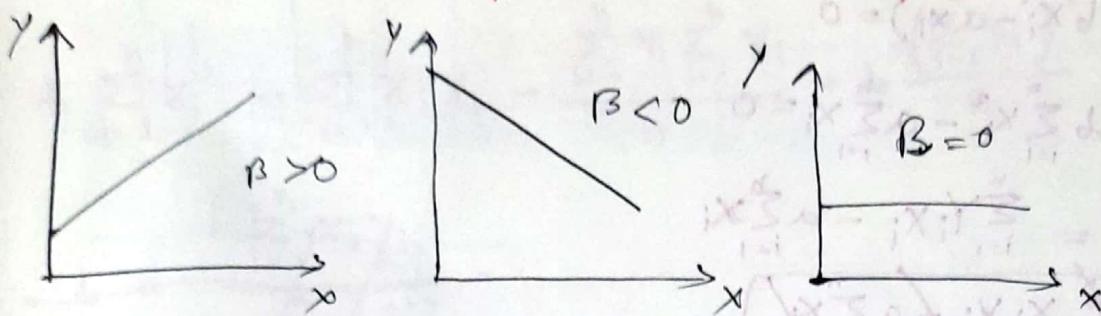
Dependent variable

$$Y = f(x)$$

Independent variable

# Simple linear dependency

$$Y = \beta X + \alpha$$



Deterministic  
Relationship  
 $Y = \beta X + \alpha$

Can predict accurate  
point when  $\beta, \alpha$  is known

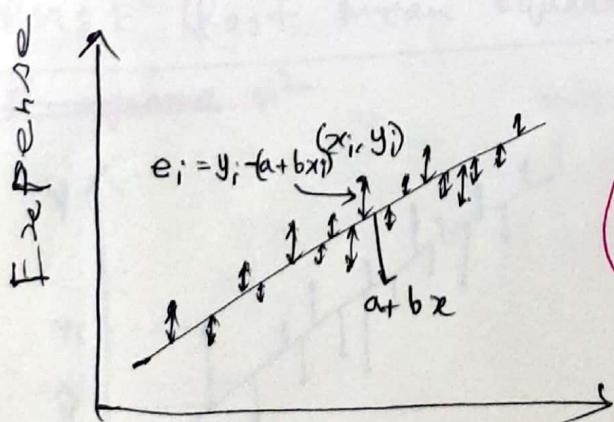
Non-deterministic  
Relationship  
 $Y = \beta X + \alpha + \epsilon$

Can not predict accurate point  
with only  $\beta$  &  $\alpha$  - slight randomness

## Simple - Linear Regression Model :

↓  
No. of independent variable 1 (single) → relationship between independent & dependent variable is linear (straight line)

Predicted value  $\rightarrow \hat{Y} = bX + a$



$$\begin{aligned} \sum (Y_i - \hat{Y})^2 \\ = \sum (Y_i - bX_i - a)^2 \end{aligned} \quad \left. \begin{array}{l} \text{error} \\ \text{calculation} \end{array} \right\}$$

Sum of square errors,

$$SSE = \sum (Y_i - bX_i - a)^2$$

$$\frac{\partial SSE}{\partial a} = -2 \sum_{i=1}^n (Y_i - bX_i - a) = 0$$

$$\Rightarrow na = \sum_{i=1}^n Y_i - b \sum_{i=1}^n X_i \quad \left[ \begin{array}{l} \therefore a = \frac{\sum_{i=1}^n Y_i}{n} - b \frac{\sum_{i=1}^n X_i}{n} \\ \therefore a = \bar{Y} - b \bar{X} \end{array} \right] \Rightarrow \sum_{i=1}^n Y_i - b \sum_{i=1}^n X_i - na = 0$$

$$\frac{\delta SSE}{\delta b} = -2 \sum_{i=1}^n (Y_i - b x_i - a) x_i = 0$$

$$\Rightarrow \sum_{i=1}^n (Y_i x_i - b x_i^2 - a x_i) = 0$$

$$\Rightarrow \sum_{i=1}^n Y_i x_i - b \sum_{i=1}^n x_i^2 - a \sum_{i=1}^n x_i = 0$$

$$\Rightarrow b \sum_{i=1}^n x_i^2 = \sum_{i=1}^n Y_i x_i - a \sum_{i=1}^n x_i$$

$$\Rightarrow b = \frac{\sum_{i=1}^n x_i Y_i - a \sum_{i=1}^n x_i}{\sum_{i=1}^n x_i^2}$$

$$\Rightarrow b \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i Y_i - (\bar{Y} - b \bar{x}) \sum_{i=1}^n x_i$$

$$\Rightarrow b \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i Y_i - \left( \frac{\sum_{i=1}^n Y_i}{n} - b \frac{\sum_{i=1}^n x_i}{n} \right) \sum_{i=1}^n x_i$$

$$b = \frac{n \sum x_i Y_i - \sum x_i \sum Y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

Ergebnis der Regressionsrechnung  
 bestimmt die Abhängigkeit zwischen abhängigen und unabhängigen Variablen  
 (mit Hilfe von Korrelationen)

$$\text{verbalisiert } \hat{Y} = (\bar{Y} - b \bar{x}) + b x$$

statische Formel:

$$(x - \bar{x})(y - \bar{y}) = 322$$

$$0 = (x - \bar{x})(y - \bar{y}) \Leftrightarrow x - \bar{x} = \frac{322}{y - \bar{y}}$$

$$0 = (x - \bar{x})(y - \bar{y}) \Leftrightarrow x - \bar{x} = \frac{322}{y - \bar{y}}$$

10. 05. 2023

Wednesday

$$b \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i - \frac{\sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n} + b \frac{(\sum_{i=1}^n x_i)^2}{n}$$

$$\cancel{b} \left[ \sum_{i=1}^n x_i y_i - \frac{\sum_{i=1}^n x_i^2}{n} \right]$$

$$b \left[ n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2 \right] = n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i$$

$$\therefore b = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

Example:

Find  $a, b$

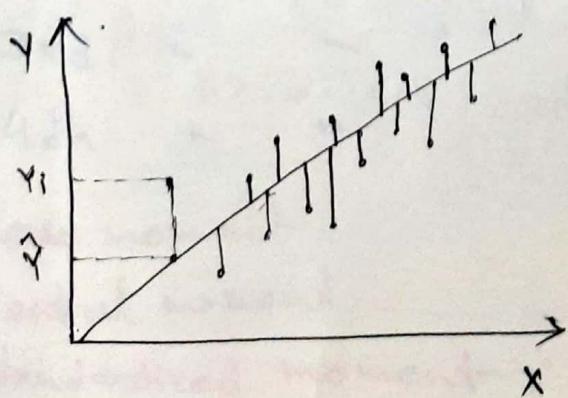
Find  $\hat{Y} = bx + a$

Find  $\hat{Y}$  for a given  $x$

When  $X$  dependent &  $Y$  independent variable,  $X = f(Y) = dY + c$

MSE (Mean square error)	$= \frac{\sum (Y_i - \bar{Y})^2}{n}$
RMS E (Root mean square error)	$= \sqrt{\frac{\sum (Y_i - \hat{Y})^2}{n}}$

R-squared  $r^2$



$$SST \text{ (Total sum of square)} = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$SSR \text{ (sum of square of regression)} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$

$$SSE \text{ (sum of square error)} = \sum_{i=1}^n (\hat{Y}_i - Y_i)^2$$

$$r^2 = \frac{SSR}{SST} = \frac{\sum_{i=1}^n (Y_i - \hat{Y})^2}{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2} = 0.5$$

can explain 50% variance of the dependent variable

- When  $r^2 \approx 0$  it doesn't hold linear regression

linear

Multiple Regression Model :

$$\hat{Y} = b_1 X_1 + b_2 X_2 + a$$

Polynomial Regression Model :

$$\hat{Y} = a + b_1 X + b_2 X^2$$

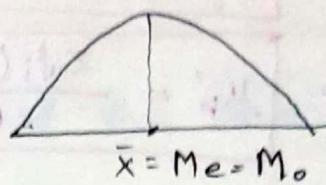
22.05.2023

Monday

## Skewness :

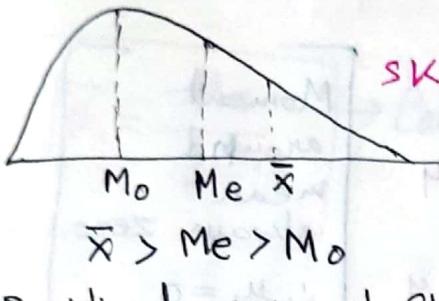
perfect

For normal distribution  $\bar{x} = M_e = M_o$



$$SK_B = 0$$

$$SK_p = 0$$

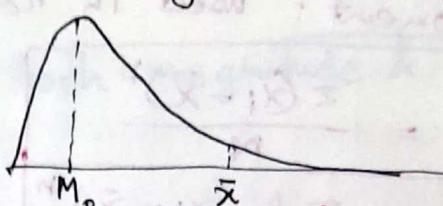


$$SK_B > 0$$

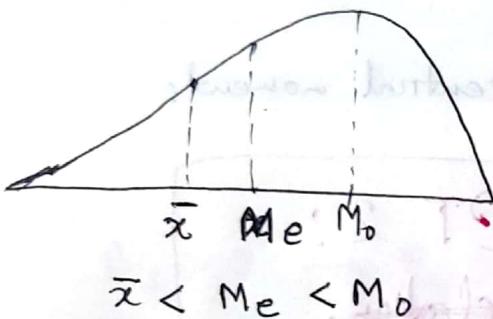
$$SK_p > 0$$

$$\bar{x} > Me > Mo$$

Positively ~~symmet~~ Skewed Distribution



Symmetrical Distribution



$$SK_B < 0$$

$$SK_p < 0$$

$$\bar{x} < Me < Mo$$

Negatively Skewed Distribution

Karl Pearson's coefficient

$$\therefore \text{Skewness, } SK_p = \frac{\text{Mean} - \text{Mode}}{\text{Standard Deviation}}$$

so that translation & scaling invariance doesn't affect skewness.

## Moments: Distance & Direction

provide unique characterization of a distribution.

Unification of all measures (central tendency, variation).

1st order moment: refers the mean of dataset

2nd  $n$   $n$ : ~~mean~~  $n$  variance  $n$   $n$

3rd  $n$   $n$ :  $n$   $n$  skewness  $n$   $n$

4th  $n$   $n$ :  $n$   $n$  kurtosis  $n$   $n$

Raw moment

Central moment

Standardized moment

### Raw moment:

ungrouped data,  $\mu'_p = \frac{\sum (x_i - A)^p}{N}$  Assumed mean  $p = 1, 2, 3, 4$

grouped data,  $\mu'_p = \frac{\sum f_i (x_i - A)^p}{N}$   $p = 1, 2, 3, 4$

with respect to origin,

$$\mu'_p = \frac{\sum x_i^p}{N}$$

$$\mu'_p = \frac{\sum f_i x_i^p}{N}$$

### Central moment: used in real life

ungrouped,  $\mu_p = \frac{\sum (x_i - \bar{x})^p}{N}$   $p = 1, 2, 3, 4$

grouped,  $\mu_p = \frac{\sum f_i (x_i - \bar{x})^p}{N}$   $p = 1, 2, 3, 4$

Moment around mean always zero.  
 $\therefore \mu_1 = 0$

### Relation between raw moment & central moment:

$$\mu_1 = 0$$

$$\mu_2 = \mu'_2 - \mu'_1^2$$

$$\mu_3 = \mu'_3 - 3\mu'_2 \mu'_1 + 2\mu'_1^3$$

$$\mu_4 = \cancel{\mu'_4}$$

Q.T 2:

Relation

Correlation,

Regression,

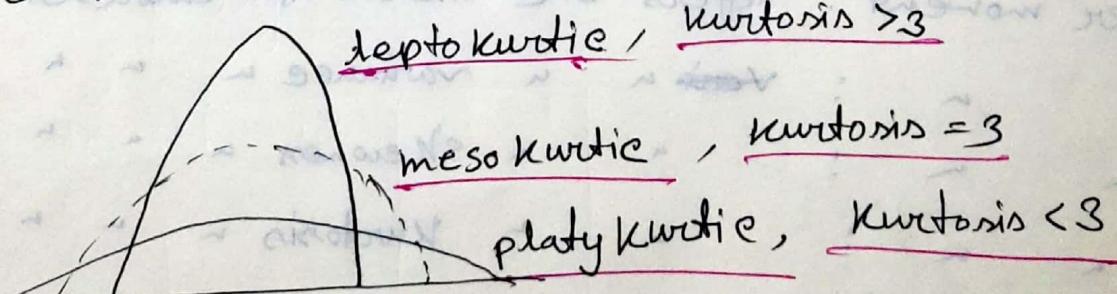
SKEWNESS,

Moment,

KURTOSIS

### Kurtosis:

measure the peakedness



Measures of skewness & kurtosis using moments.

Relative measure of skewness,  $\beta_1 = \frac{\mu_3}{\mu_2^{3/2}}$

always positive

$$\mu_3^2 / \mu_2^3$$

variance

$\beta_1$  can determine the magnitude of the skewness but not the direction

$$\gamma_1 = \sqrt{\beta_1} = \frac{\mu_3}{\mu_2^{3/2}}$$

Can determine both magnitude & direction

Relative measure of kurtosis,

$$\beta_2 = \frac{\mu_4}{\mu_2^2}$$

$$\gamma_2 = \beta_2 - 3$$

excess kurtosis

$\beta_2 > 3 \rightarrow \gamma_2 > 0 \rightarrow$  leptokurtic

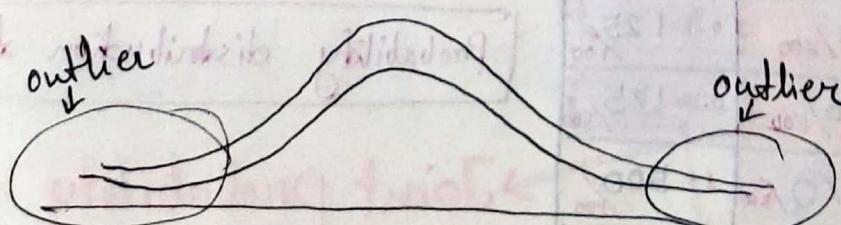
$\beta_2 = 3 \rightarrow \gamma = 0 \rightarrow$  mesokurtic

$\beta_2 < 3 \rightarrow \gamma < 0 \rightarrow$  platykurtic

To measure kurtosis we used to concentrate on peak of the data but nowadays we concentrate more on tail or outlier of the data.

amount of outlier  $\uparrow$  kurtosis  $\uparrow$

amount of outlier  $\downarrow$  kurtosis  $\downarrow$



leptokurtic	mesokurtic	platykurtic
thin tails	medium tails	wide tails

24.05.23  
Wednesday

## Descriptive Statistics

### Univariate

- | Mean
- | Median
- | Mode
- | Variance
- | Standard deviation
- | Normal distribution / z-score
- | Skewness
- | Kurtosis
- | Moments

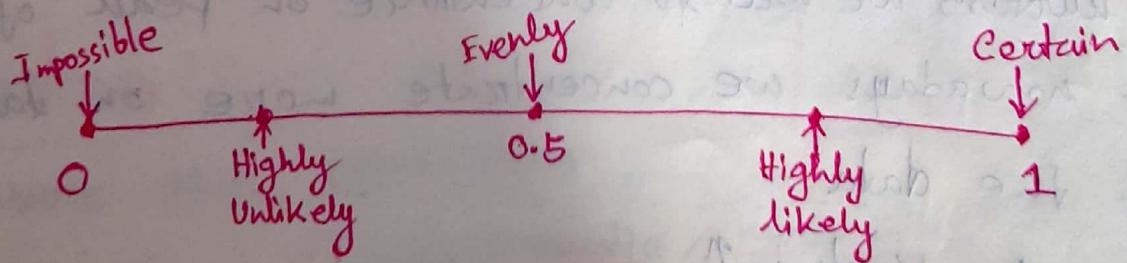
### Multivariate

- | Covariance
- | Correlation
- | Regression

## Probability :

- Joint probability
- Marginal probability
- Conditional probability
- Independence

## Event probability :



#	Male	Female	Total
Friends	$0.16 = \frac{80}{500}$	$0.24 = \frac{120}{500}$	$0.4 = \frac{200}{500}$
Big Bang Theory	$0.2 = \frac{100}{500}$	$0.05 = \frac{25}{500}$	$0.25 = \frac{125}{500}$
Others	$0.1 = \frac{50}{500}$	$0.25 = \frac{125}{500}$	$0.35 = \frac{175}{500}$
Total	$0.46 = \frac{230}{500}$	$0.54 = \frac{270}{500}$	$1 = \frac{500}{500}$

Probability distribution table

Joint probability  
Marginal probability distribution {P(A)}

Distribution {P(AnB)}

Marginal probability: single term is represented  
on ~~single~~ probability  
simple

\*  $P(G = M \text{ OR } F = \text{Big Bang Theory}) = 0.51 = 0.16 + 0.2 + 0.1 + 0.05$

\*  $P(A \cup B) = P(A) + P(B) - P(A \cap B) = 0.46 + 0.25 - 0.2 = 0.51$

\* We know that subscriber is female then what is the probability of first choice to be Big Bang theory?

$$P(F = "BBT" | G = 'F') = \frac{0.05}{0.54} = 0.09$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

an event has already occurred and : Conditional probability  
another will happen w.r.t. it.

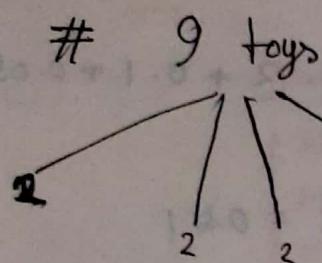
Independent event:  $P(A|B) = P(A)$

No relation between events A & B

If two events are  $\Rightarrow$  independent  $P(A \cap B) = P(A) \cdot P(B)$

29.05.2023

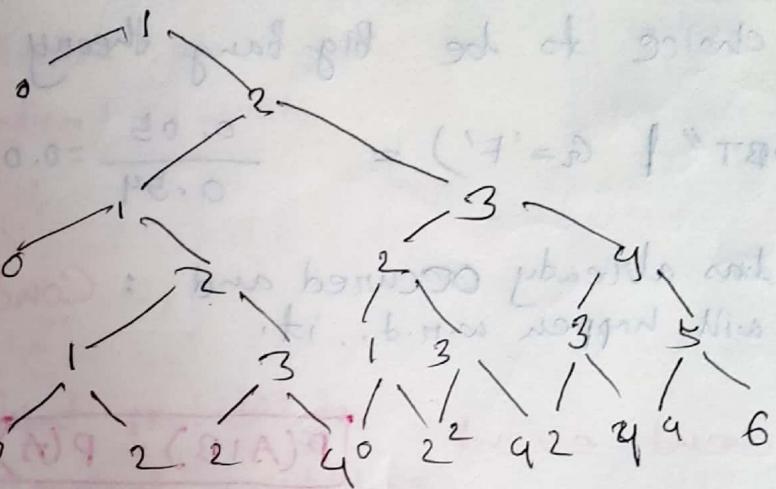
Monday



$$9c_3 \times 6c_2 \times 4c_2 \times 2c_2 = 7560 = \frac{9!}{2!2!2!3!} = \frac{123456789}{22222222}$$

$$\frac{2 \cdot 3 \cdot 3^4}{(3!)^3}$$

$$\frac{(3!)^3}{(3!)^3} = \frac{1}{(3!)^3}$$



a) 14 b) 2

c) 4

Poker:

Royal flush: 10 J Q K A  $\rightarrow$  same suit (4 possible)

probability :  $\frac{4}{52c_5}$

Straight flush: 5 consecutive card  $\rightarrow$  same suit

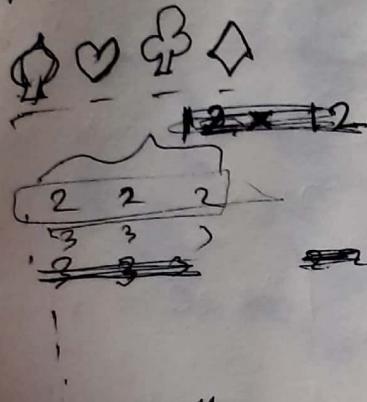
$$\frac{9 \times 4}{52c_5}$$

$\left\{ \begin{array}{l} 2 3 4 5 6 \\ 3 4 5 6 7 \\ 4 5 6 7 8 \\ 5 6 7 8 9 \\ 6 7 8 9 10 \end{array} \right.$	$\left. \begin{array}{l} 10 \\ A \end{array} \right\} \rightarrow \text{Royal flush}$
$\left\{ \begin{array}{l} 2 3 4 5 6 \\ 3 4 5 6 7 \\ 4 5 6 7 8 \\ 5 6 7 8 9 \\ 6 7 8 9 10 \end{array} \right.$	$\left. \begin{array}{l} 10 \\ A \end{array} \right\} \rightarrow \text{Royal flush}$
$\left\{ \begin{array}{l} 2 3 4 5 6 \\ 3 4 5 6 7 \\ 4 5 6 7 8 \\ 5 6 7 8 9 \\ 6 7 8 9 10 \end{array} \right.$	$\left. \begin{array}{l} 10 \\ A \end{array} \right\} \rightarrow \text{Royal flush}$
$\left\{ \begin{array}{l} 2 3 4 5 6 \\ 3 4 5 6 7 \\ 4 5 6 7 8 \\ 5 6 7 8 9 \\ 6 7 8 9 10 \end{array} \right.$	$\left. \begin{array}{l} 10 \\ A \end{array} \right\} \rightarrow \text{Royal flush}$

Four of a Kind : 4 card same number , 1 card different

$$\frac{13 \times (52-4)}{52C_5}$$

Full house : 3 of a kind (3 card same number), 1 pair  
(2 card same)

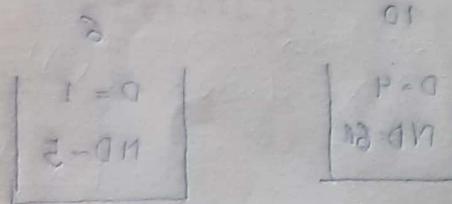


$$\frac{13 \times 4C_3 \times 12 \times 4C_2}{52C_5}$$

$$(11)9 \cdot (8)9 = (88)9$$

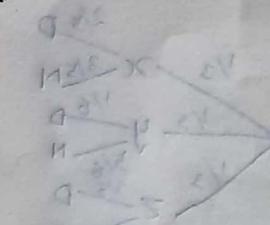
Flush : 5 card same suit

$$\frac{13C_5 \times 4}{52C_5} = \frac{4 + (9 \times 4)}{52 - 11}$$



~~Straight~~ Straight : 5 consecutive card → suit doesn't matter

$$\frac{(10 \times 4^5) - (40)}{52C_5} = (41)9$$



$$\frac{(11)9 \cdot (9)9}{(11)9 \cdot (10)9} = \frac{(10)9}{(11)9} = (11)9$$

$$\frac{(11)9 \cdot (10)9}{(11)9 \cdot (11)9} = (10)9$$

8 P. Diagram

31.05.2023  
Wednesday

## Conditional Probability

### Bayes Theorem

#### Naive Bayes Theorem

Conditional Probability:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

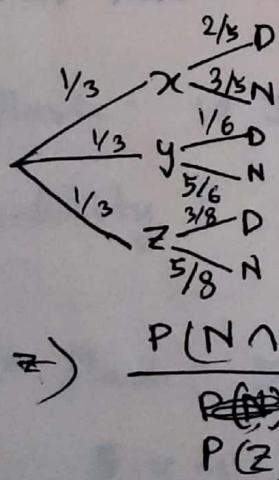
$$\boxed{P(A \cap B) = P(B) \cdot P(A|B)}$$

dependent

$$\begin{array}{|c|} \hline D=4 \\ ND=60 \\ \hline \end{array}$$

$$\begin{array}{|c|} \hline D=1 \\ ND=5 \\ \hline \end{array}$$

$$\begin{array}{|c|} \hline D=3 \\ ND=5 \\ \hline \end{array}$$



$$P(D) = \frac{1}{3} \cdot \frac{2}{5} + \frac{1}{3} \cdot \frac{1}{6} + \frac{1}{3} \cdot \frac{3}{8}$$

$$P(ND) = 1 - P(D)$$

$$P(N|Z) = \frac{P(N \cap Z)}{P(Z)} = \frac{P(Z) \cdot P(N|Z)}{P(X) \cdot P(N|X) + P(Y) \cdot P(N|Y) + P(Z) \cdot P(N|Z)}$$

#### Bayes Theorem

$$P(Z|N) = \frac{P(N \cap Z)}{P(N)} = \frac{P(Z) \cdot P(N|Z)}{P(X) \cdot P(N|X) + P(Y) \cdot P(N|Y) + P(Z) \cdot P(N|Z)}$$

# Example 4.8

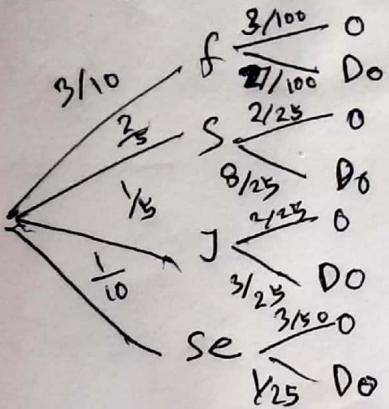
1.

$$P(F) = 30\%$$

$$P(S) = 40\%$$

$$P(J) = 20\%$$

$$P(Se) = 10\%$$



$$P(F|O) = 10\%$$

$$P(S|O) = 20\%$$

$$P(J|O) = 40\%$$

$$P(Se|O) = 60\%$$

$$P(O) = \frac{3}{100} + \frac{2}{25} + \frac{2}{25} + \frac{3}{50}$$

$$= 0.25$$

$$P(J|O) = \frac{P(J \cap O)}{P(O)} = \frac{P(J) \cdot P(O|J)}{P(O)}$$

$$= \frac{2/25}{0.25}$$

Independent :

$$P(A \cap B) = P(A) \cdot P(B)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

09.07.2023

Sunday

## Chapter 2, 3, 4, 5

$$P(c_i | x) = \frac{P(x | c_i) * P(c_i)}{P(x)}$$

$c_i$  = Class we belong to (Yes/ No) class

$x$  = feature vector

$P(\text{Yes} | x)$

$P(\text{No} | x)$

Need to find

$$P(x | c_i) * P(c_i)$$

$$P(x | c_i) = \prod P(x_1 | c_i) P(x_2 | c_i) \dots P(x_n | c_i)$$

Example:

$$X = \underbrace{\text{age} \leq 30}_{x_1}, \underbrace{\text{income} = \text{medium}}_{x_2}, \underbrace{\text{student} = \text{yes}}_{x_3}, \underbrace{\text{credit-rating} = \text{fair}}_{x_4}$$

find each  $P(X_n | C_i)$  from dataset, then multiply  
for both Yes & No class

the answer of  $P(x_i | C_i) * P(C_i)$  shows  $x$  belongs to which class. (it belongs to the class which has greater  $P(x_i | C_i) * P(C_i)$  value)

### Random variable: (Important) (Chapter 5) (Book)

Always map to the value of the outcome of a Random Process.  
Outcome cannot be predicted.

represented with Capital ( $X$ ) always.  
letter

$$X = \begin{cases} 0, & \text{if outcome} \rightarrow \text{tail} \\ 1, & \text{if outcome} \rightarrow \text{head} \end{cases}$$

- \* Discrete random variable (Countable)(finite)  
Counted
- \* Continuous random variable (cannot be counted exactly)  
(e.g., predicting the amount of rain tomorrow)

Example 5.3:  $X = \max(a, b) = \{1, 2, 3, 4, 5, 6\}$  (all possible values)  
we can count it, thus it is discrete

$$Y = \sum(a, b) = 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12$$

$X$ :

$$f(1) = (1, 1) \rightarrow 1/36$$

$$f(2) = (1, 2), (2, 1) \rightarrow 2/36$$

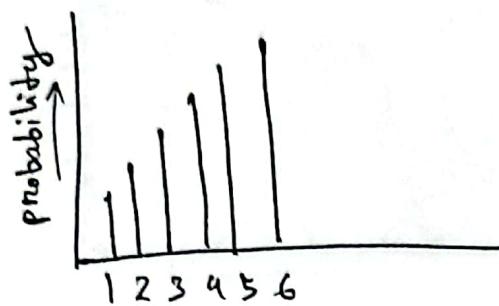
$$f(3) = (1, 3), (2, 2), (3, 1) \rightarrow 3/36$$

$$f(4) = 4/36, f(5) = 3/36, f(6) = 2/36$$

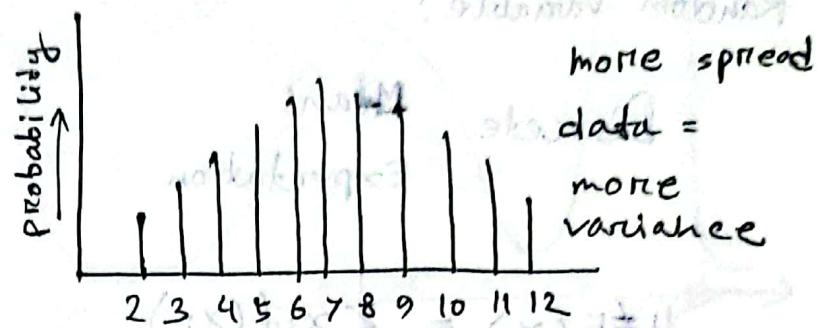
$Y$ : same way find for  $Y$

Distribution:

Distribution of  $X$ :



Distribution of  $Y$ :



Expectation of a finite random variable:  
(Mean)

$$E = E(X) = x_1 f(x_1) + x_2 f(x_2) + \dots + x_n f(x_n) = \sum x_i f(x_i)$$

$$\text{or, } E = E(X) = x_1 p_1 + x_2 p_2 + \dots + x_n p_n = \sum x_i p_i,$$

$$\text{where, } p_i = f(x_i) = P(X = x_i)$$

Example 5.8:

x <sub>i</sub>	0	1	2	3	4	5	6
p <sub>i</sub>	1/64	6/64	15/64	30/64	15/64	6/64	1/64

Probability of getting head when tossing 6 times.

(use  $nCr$ )

Example 5.9: a dice is tossed. if value 2, 3, 5  $\rightarrow$  wins  
if value 1, 4, 6  $\rightarrow$  losses

x	2	3	5	-1	-4	-6
f(x)	1/6	1/6	1/6	1/6	1/6	1/6

$$E(x) = 2(1/6) + 3(1/6) + 5(1/6) - 1(1/6) - 4(1/6) - 6(1/6) = -1/6$$

not a favorable game as  $E(x)$  is negative. the game would have been fair if  $E(x)$  is 0

10.07.2023

Monday

Random variable:

Discrete, Mean:

Expectation

$$\mu = E(x) = \sum x_i f(x_i)$$

Probability distribution function

$$f(x_i) \geq 0 ; \sum f(x_i) = 1$$

$$\text{Var}(x) = (x_1 - \mu)^2 f(x_1) + (x_2 - \mu)^2 f(x_2) + \dots + (x_n - \mu)^2 f(x_n)$$

$$= \sum (x_i - \mu)^2 f(x_i)$$

$$= E(x^2) - \mu^2$$

$$= \sum x_i^2 f(x_i) - \mu^2$$

$$\sigma = \sqrt{\text{Var}(x)}$$

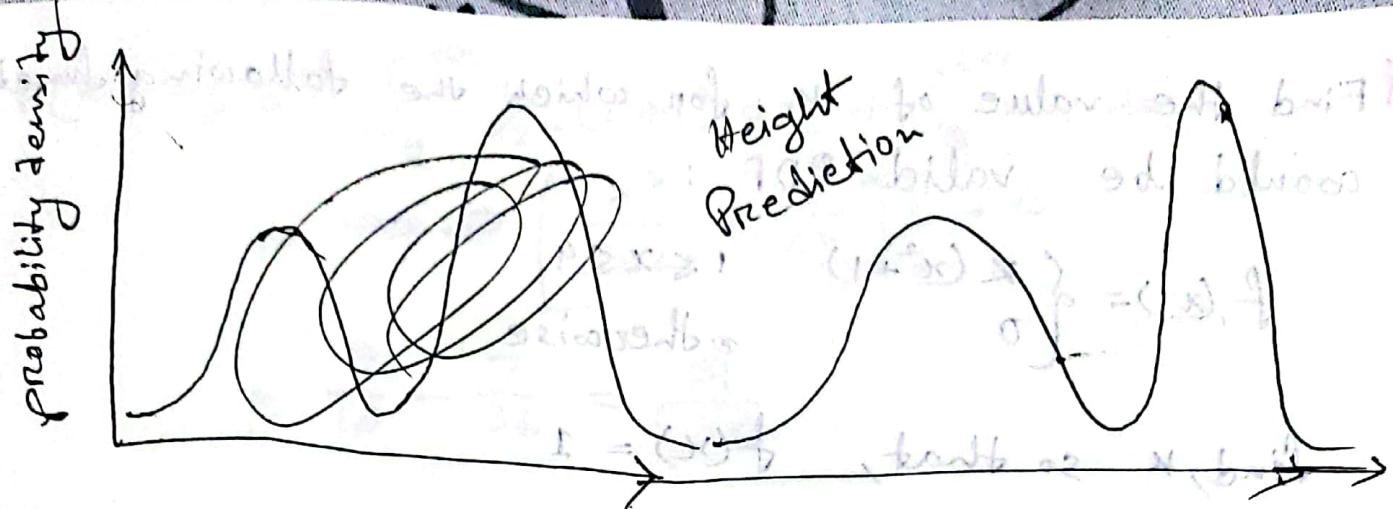
Example 5.11

Standardized Random variable (Later)

Joint random variable (Skip)

Continuous Random Variable:

Here,  $f(x_i) \rightarrow$  Probability Density Function (PDF)



Here we work with intervals instead of exact point.

$$P(100 < x < 150) = \int_{100}^{150} f(x) dx$$

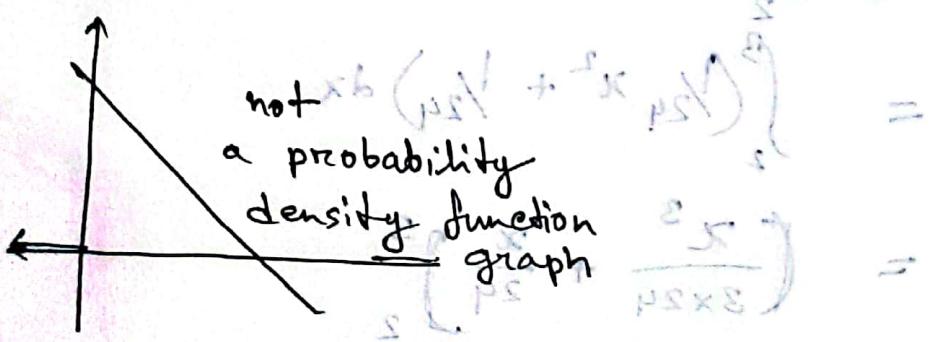
$$P(a < x < b) = \int_a^b f(x) dx$$

The total area under the curve,  $\int_{-\infty}^{\infty} f(x) dx = 1$

Real life graph function:

$f(x)$  = very complicated

thus we use more theoretically made up graph.



To find whether a graph is PDF or not:

i) Is the area  $> 0$ ?

ii) Is the area / probability = 1?

\* Find the value of  $\kappa$  for which the following function could be valid PDF:

$$f(x) = \begin{cases} \kappa(x^2+1) & 1 \leq x \leq 4 \\ 0 & \text{otherwise} \end{cases}$$

Find,  $\kappa$  so that,  $f(x) = 1$

Ans.:  $\kappa = 1/24$

\*  $f(x) = \begin{cases} 1/24(x^2+1) & ; 1 \leq x \leq 4 \\ 0 & ; \text{otherwise} \end{cases}$

i)  $P(x=2)$

ii)  $P(2 \leq x \leq 3)$

iii)  $P(x \geq 2)$

i)  $P(x=2) = 0$

ii)  $P(2 \leq x \leq 3) = \int_{2}^{3} 1/24(x^2+1) dx$

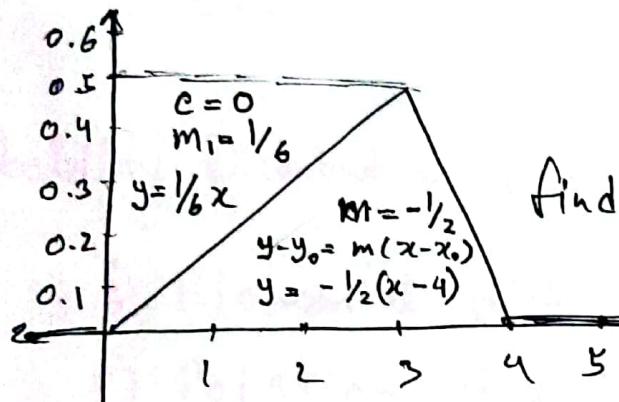
$$= \int_{2}^{3} (1/24x^2 + 1/24) dx$$

$$= \left[ \frac{x^3}{3 \times 24} + \frac{x}{24} \right]_{2}^{3}$$

$$\Rightarrow \frac{11}{26}$$

$$\begin{aligned}
 \text{- III) } P(x \geq 2) &= \int_2^{4\cancel{0}} \frac{1}{24}(x^2+1) dx + \int_{4\cancel{0}}^{\infty} \dots dx \\
 &= \left[ \frac{x^3}{3 \cdot 24} + \frac{x}{24} \right]_2^4 + 0 \\
 &= \frac{19}{18} - \frac{7(5)}{36} = \frac{31}{36}
 \end{aligned}$$

#



$$f(x) = \begin{cases} \frac{1}{6}x & ; 0 \leq x \leq 4 \\ -\frac{1}{2}(x-4) & ; 4 < x < 5 \end{cases}$$

Cumulative Density Function (CDF) → next class

Cumulative Density Function

$$F(x_0) = P(x \leq x_0)$$

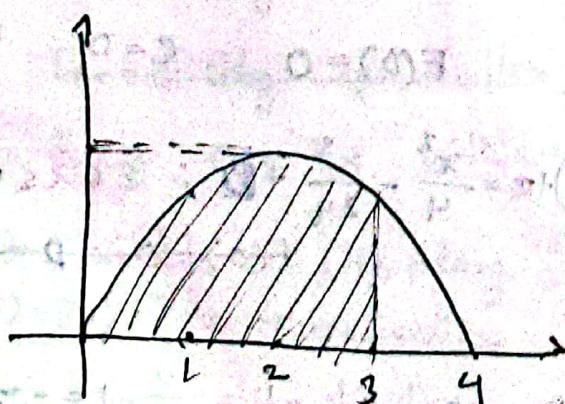
summation of density function

17.07.2023

Monday

$$\text{Discrete: } F(x_0) = \sum_{x_i < x_0} f(x_i)$$

$$\text{Continuous: } F(x_0) = \int_0^{x_0} f(x) dx$$



Probability Density Function (PDF)  
Cumulative Distribution Function (CDF)  
identify first when doing math



$$P(2 < x < 3) \rightarrow F(3) - F(2)$$

$$P(x > 2) \rightarrow 1 - F(2)$$

$$\begin{cases} 0 & x < 0 \\ \frac{1}{32}(6x^2 - x^3) & 0 \leq x \leq 4 \\ 1 & x > 4 \end{cases}$$

Integrate

PDF  $\rightarrow$  CDF

Differentiate

CDF  $\rightarrow$  PDF

Conversion

$$\begin{cases} \frac{1}{8}x^2 & ; 0 \leq x \leq 2 \\ \frac{1}{8}x(4-x) & ; 2 < x \leq 4 \\ 0 & ; \text{otherwise} \end{cases}$$

$$\begin{cases} 0 & ; x < 0 \\ \frac{1}{24}x^3 & ; 0 \leq x \leq 2 \\ \frac{x^2}{4} - \frac{x^3}{24} - \frac{1}{3} & ; 2 < x \leq 4 \\ 1 & ; x > 4 \end{cases}$$

$$F(x) = \int \frac{1}{8}x^2 dx = \frac{1}{24}x^3 + C; 0 \leq x \leq 2 \quad F(0) = 0 \therefore C = 0$$

$$F(x) = \int \frac{1}{8}(4x - x^2) dx = \frac{1}{8}\left(4x^2 - \frac{x^3}{3}\right) + D = \frac{x^3}{4} - \frac{x^3}{24} + D; 2 \leq x \leq 4$$

$$F(2) = \frac{1}{24}2^3 = \frac{8}{24} = \frac{1}{3}$$

$$F(0) = 0 \therefore D = 0$$

$$\therefore \frac{1}{3} = \frac{2^3}{4} - \frac{2^3}{24} + D \quad \therefore D = \frac{1}{3} + \frac{1}{3} - 1 = \frac{2}{3} - 1 = -\frac{1}{3}$$

## Mean & Variance of Continuous Random Variable

$$\mu = E(x) = \sum x_i f(x_i)$$

$$\sigma^2 = E(x^2) - \mu^2$$

$$\mu = \int_{-\infty}^{\infty} x f(x) dx$$

$$\sigma^2 = \int_{-\infty}^{\infty} x^2 f(x) dx - \mu^2$$

Discrete random variable

Continuous random variable

## Probability Distribution of Discrete Random Variable

i) Binomial distribution

ii) Poisson distribution

### Binomial Distribution:

Can be applied to problems which have specific two outcomes  
 Success      Failure

### Bernoulli trial

Process of occurrence  $\rightarrow$  Bernoulli

Each trial is independent.

Probability of success ( $P$ )  $\rightarrow$  remain constant

# Probability of an item 3 items are picked are random,  
 they can be either defective or non-defective.

$$P(D) = 25\% = \frac{1}{4} \quad | \quad 3 \rightarrow \text{Display}$$

$$P(N) = 75\% = \frac{3}{4} \quad | \quad 2 \rightarrow \text{defective}$$

success ( $P$ )

failure.  $(1-P)/(n)$

$$P(\text{2 defective}) = 3 \times \left(\frac{1}{4}\right)^2 \times \frac{3}{4}$$

$$S = \{NNN, NDN, NND, DNN, \underline{N DD}, \underline{DND}, \underline{DDN}, DDD\} = \frac{9}{64}$$

$$P(x \text{ success}) = {}^n C_x \times P^x \times (1-P)^{n-x}$$

$P$  = success

$(1-P)$  = failure

# 10 → display →  $n$   
 5 → non defective →  $n-x$

$$P(\text{success}) = {}^{10} C_5 \times (1/4)^5 \times (3/4)^{10-5}$$

10 → display

at least 5 → defective

$$\begin{aligned} P(x \text{ success}) &= P(5) + P(6) + P(7) + P(8) + P(9) + P(10) \\ &= b(5, 10, 1/4) + b(6, 10, 1/4) + b(7, 10, 1/4) + \dots \end{aligned}$$

$$= \sum_{x=5}^{10} \left\{ {}^{10} C_x \times (1/4)^x \times (3/4)^{10-x} \right\}$$

Use calculator

$$= 1 - \sum_{x=0}^4 \left\{ {}^{10} C_x \times (1/4)^x \times (3/4)^{10-x} \right\}$$

10 → display

5 - 7 → defective

$$\sum_{x=0}^7 - \sum_{x=0}^4$$

18.07.2023

Tuesday

### Discrete

Binomial distribution

Poisson distribution

### Continuous

Uniform distribution

Normal distribution

$$\# a) \sum_{x=1}^{20} {}^{20}C_x \left(\frac{3}{100}\right)^x \left(\frac{97}{100}\right)^{20-x}$$

$$\text{or } = 1 - {}^{20}C_0 \left(\frac{3}{100}\right)^0 \times \left(\frac{97}{100}\right)^{20}$$

$$= 0.4562$$

$$b) {}^{10}C_3 \left(0.4562\right)^3 \left(1-0.4562\right)^7 = 0.16023$$

Mean & variance in binomial distribution

$$X = I_1 + I_2 + \dots + I_n$$

$$\mu = E(X) = E(I_1) + E(I_2) + \dots + E(I_n) = \underbrace{p + p + \dots + p}_{n \text{ terms}} = np$$

$I_j \rightarrow$  <sup>0 or 1</sup>  $\rightarrow$  mean  $p = \bar{x}$  to optimising to probability

$$I_j \rightarrow E(I_j) = (0)q + (1)p = p$$

$$\sigma^2_{I_1} = \sum x^2 f(x_i) - \mu^2$$

$$= 0^2 \times q + 1^2 \times p - p^2$$

$$= p - p^2 = p(1-p) = pq$$

$$\therefore \sigma^2_X = \sigma^2_{I_1} + \sigma^2_{I_2} + \dots + \sigma^2_{I_n}$$

$$= \underbrace{pq + pq + \dots + pq}_{n \text{ terms}}$$

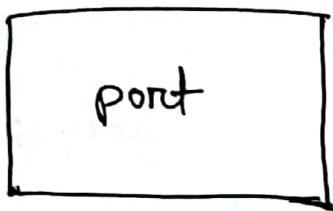
$$= npq$$

Poisson distribution: एकीं Random event or average rate द्वारा आते, एथन निम्नलिखित असंघर्षक रूप से घटना होती है probability के बराबर यह Poisson distribution कहा जाता है।

$$P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

#

Average = 10 oil tankers



Maximum capacity  
= 15 tankers.

What is the probability on a given day, a tanker has to be turned away?

$$\therefore P(X > 15) \rightarrow ?$$

$$= 1 - \sum_{x=0}^{15} \frac{\lambda^x e^{-\lambda}}{x!} = 1 - \sum_{x=0}^{15} \frac{10^x e^{-10}}{x!}$$

Probability of possibility of 3-7 tanks being available.

$$= \sum_{x=0}^7 \frac{10^x e^{-10}}{x!} - \sum_{x=0}^2 \frac{10^x e^{-10}}{x!}$$

$$\text{OR} = \sum_{x=3}^7 \frac{10^x e^{-10}}{x!}$$

# probability of accident to occur (average) = 0.005

In the next 400 days exactly once the accident will occur → what is the probability of it?

$$\lambda = 400 \times 0.005 = 2$$

$$\therefore \frac{2^2 e^{-2}}{2!} = \frac{2}{e^2} (\text{Answer})$$

# Relation between binomial & poisson distribution

trial number of binomial distribution,  $n \rightarrow \infty$

probability of success,  $p \rightarrow 0$

$np \rightarrow \text{constant}$

$$nC_x p^x q^{n-x} \underset{n \rightarrow \infty}{\approx} \frac{\lambda^x e^{-\lambda}}{x!}$$

derivation

19.02.2023

Wednesday

$nC_x p^x (1-p)^{n-x}$  expanding!

$$= \frac{n!}{x!(n-x)!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \rightarrow 1 \cdot (1 - \frac{1}{n})(1 - \frac{2}{n}) \cdots (1 - \frac{x+1}{n})$$

$$= \frac{\lambda^x}{x!} \left( \frac{n!}{(n-x)!} \cdot \frac{1}{n^x} \right) \left(1 - \frac{\lambda}{n}\right)^{n-x} e^x = \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n$$

$$= \frac{\lambda^x}{x!} \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{x+1}{n}\right) \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x}$$

$$= \frac{\lambda^x}{x!} \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{x+1}{n}\right) \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x}$$

$$\lim_{n \rightarrow \infty} nC_x p^x (1-p)^{n-x}$$

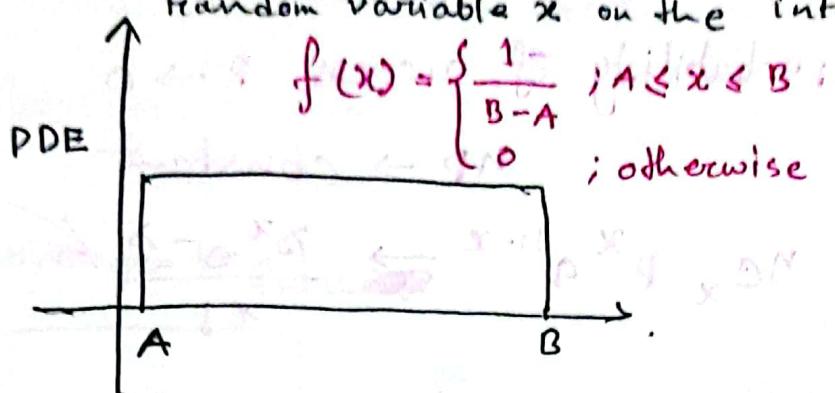
$$= \frac{\lambda^x}{x!} \cdot \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{x+1}{n}\right) x \lim_{n \rightarrow \infty} \left(1 + \frac{-\lambda}{n}\right)^n \lim_{n \rightarrow \infty} \left(1 + \frac{-\lambda}{n}\right)^{-x}$$

$$= \frac{\lambda^x}{x!} \cdot 1 \cdot e^{-\lambda} \cdot 1 = \frac{\lambda^x e^{-\lambda}}{x!}$$

Probability & statistics WAL POLE → distribution part.

## Continuous Probability Distribution:

**Uniform Distribution:** the continuous function of the continuous random variable  $x$  on the interval  $[A, B]$  is,



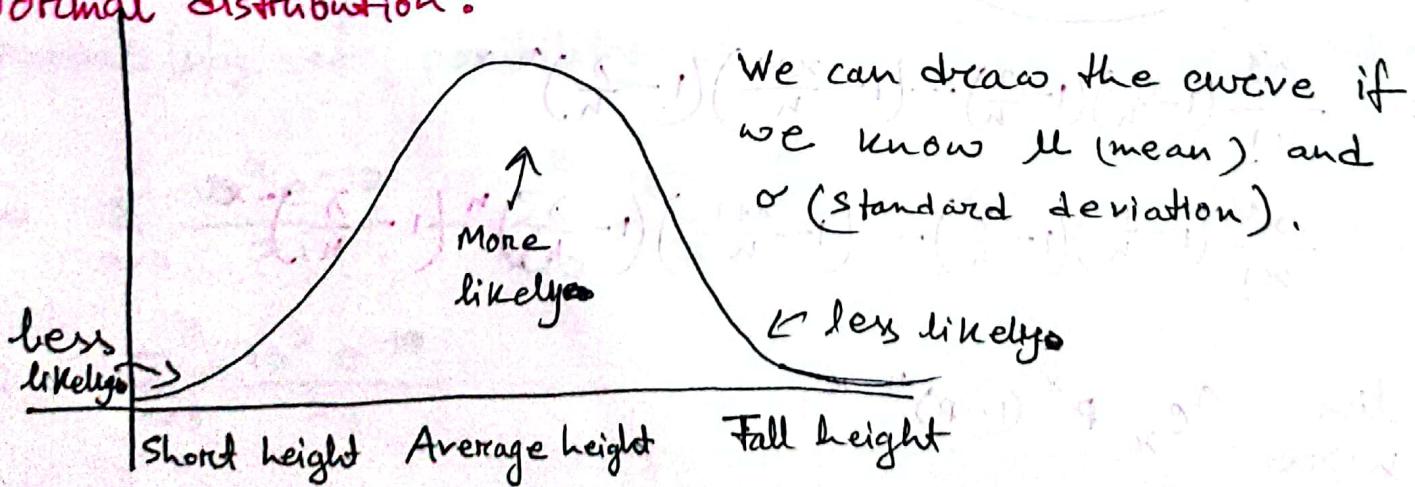
# 9)  $f(x) = \begin{cases} 1/4 & ; 0 \leq x \leq 4, \\ 0 & ; \text{elsewhere} \end{cases}$

b) Probability that any given conference lasts at least 3 hours:

$$P[x \geq 3] = \int_3^4 \frac{1}{4} dx = 1/4$$

Real-life uniform distribution is rare!

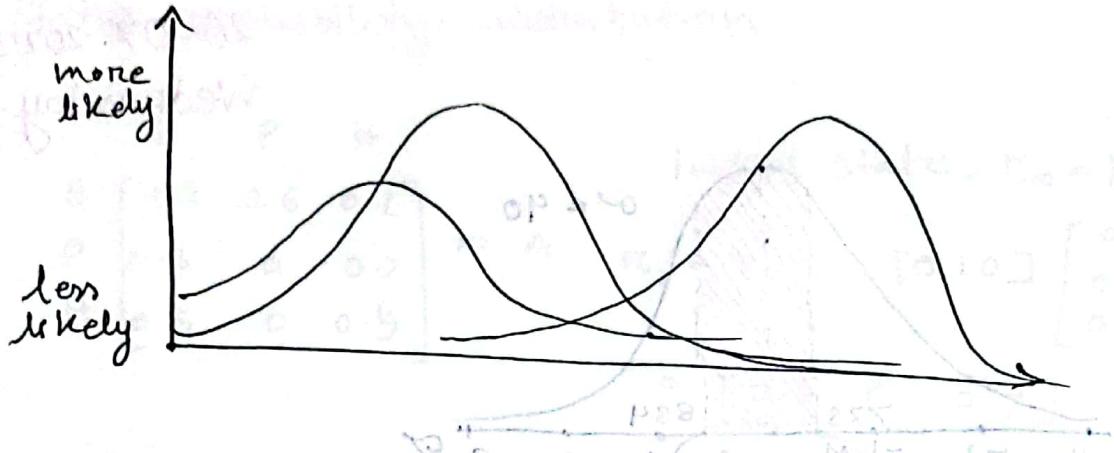
Normal distribution:



Maximum real life distribution is normal distribution due to "Central Limit Theorem"

Normal distribution function / Gaussian function / Bell shape curve:

$$n(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \quad -\infty < x < \infty$$



Standardized normal distribution:  $\mu = 0 \quad \sigma = 1$

Convert using Z-score:  $= \frac{x-\mu}{\sigma}$

Theorem

Example 6.2:

$$E(x) = E(x - \mu + \mu)$$

$$= E(x - \mu) + \mu$$

$$E(x - \mu) = \int_{-\infty}^{\infty} \frac{x - \mu}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2} dx$$

$$Z = \frac{x - \mu}{\sigma}; \quad dx = \sigma dz \rightarrow E(x - \mu) = 0$$

$$\therefore E(x) = \mu$$

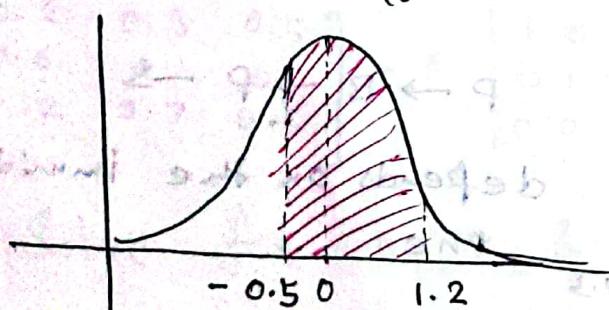
#

$$\mu = 50$$

$$\sigma = 10$$

$$Z_1 = \frac{45 - 50}{10} = -0.5$$

$$Z_2 = \frac{62 - 50}{10} = 1.2$$



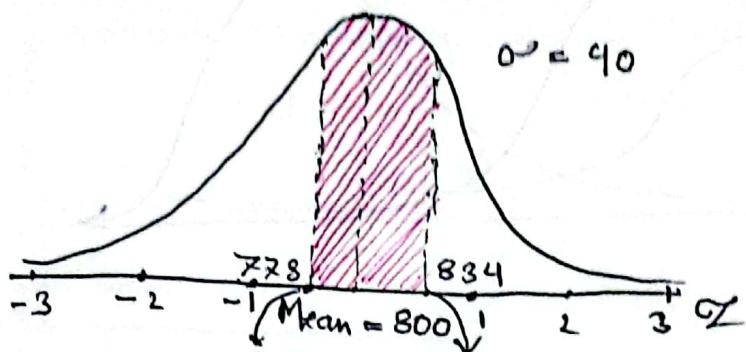
$$P(45 < x < 62) = P(-0.5 < Z < 1.2) = P(Z < 1.2) - P(Z < -0.5)$$

$$= 0.8849 - 0.3085 = 0.5764$$

Using calculator: we can find these values from Z-table

26.07.2023  
Wednesday

Example:



$$x_1 = 778 \quad x_2 = 834$$

$$z_1 = \frac{778 - 800}{40} = -0.55 \quad z_2 = \frac{834 - 800}{40} = 0.85$$

$$P(0.85) - P(-0.55) = \text{pink area}$$

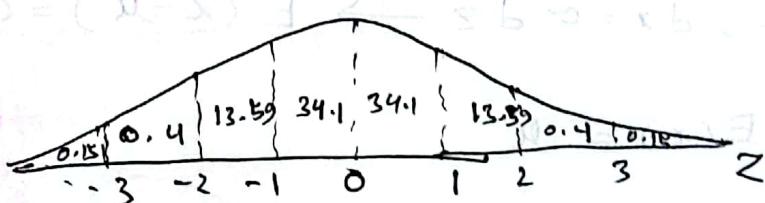
Example: Calculate  $z$  from given area (Probability):

$$\mu = 40$$

$$\sigma = 6$$

a) area 95%.

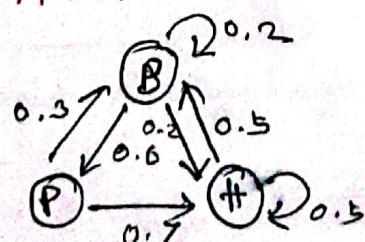
to the left



b) area 15%.  ~~$= 0.2 < z < -0.1$~~

to the right

## MARCOV CHAIN



$$P \rightarrow B \rightarrow P \rightarrow$$

depends on the immediate previous one.

After 10 steps ~~H~~  $P(B)$   $P(P)$   $P(H)$

$$\frac{4}{10}, \quad \frac{2}{10}, \quad \frac{4}{10}$$

$$B \rightarrow P \rightarrow B \rightarrow H \rightarrow B \rightarrow H \rightarrow H \rightarrow B \rightarrow H \rightarrow P$$

After  $\infty$  steps  $\rightarrow P(B) \quad P(P) \quad P(H) \rightarrow$  goes static

$$\frac{0.25}{0.45}, \quad \frac{0.3}{0.45}, \quad \frac{0.4}{0.45}$$

## Static probability distribution

$$\begin{matrix} & B & P & H \end{matrix}$$

$$B \begin{bmatrix} 0.2 & 0.6 & 0.2 \\ 0.3 & 0 & 0.7 \\ 0.5 & 0 & 0.5 \end{bmatrix} = A$$

Initial state,  $\pi_0 = [0 \ 1 \ 0]$

$$\pi_1 = \pi_0 A = [0 \ 1 \ 0] \begin{bmatrix} 0.2 & 0.6 & 0.2 \\ 0.3 & 0 & 0.7 \\ 0.5 & 0 & 0.5 \end{bmatrix} = [0.3 \ 0 \ 0.7]$$

$$\pi_2 = \pi_1 A = [0.3 \ 0 \ 0.7] \begin{bmatrix} 0.2 & 0.6 & 0.2 \\ 0.3 & 0 & 0.7 \\ 0.5 & 0 & 0.5 \end{bmatrix}$$

⋮

$$\pi_n = \pi_{n-1} A$$

→ find using linear algebra as it is not possible to pursue brute-force way in a short time.

when static →  
or stationary

Eigen value, Eigen vector :  $A v = \lambda v$

$$A \pi = 1 \cdot \pi$$

$$\vec{v} \leftarrow \pi = [\pi_1 \ \pi_2 \ \pi_3]$$

$$\pi_1 + \pi_2 + \pi_3 = 1$$

$$\det(A - \lambda I)$$

$$(A - \lambda I) \cdot \vec{v} = 0 \Rightarrow (A^T - \lambda I)^T \vec{v} = 0$$

$$\Leftrightarrow \left( \begin{bmatrix} 0.2 & 0.3 & 0.5 \\ 0.6 & 0 & 0 \\ 0.2 & 0.7 & 0.5 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \right) \cdot \begin{bmatrix} \pi_1 \\ \pi_2 \\ \pi_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Verification

⇒ Solve for  $\pi_1, \pi_2$  &  $\pi_3$ .