

# Statics

# Statistics

Statistics is the science of collecting, organizing and analyzing data.

Data:- "fact or pieces of information"

## Types of Statistics:

### ① Descriptive statistics

→ It consists of organizing and summarizing data.

#### Sub Types!:

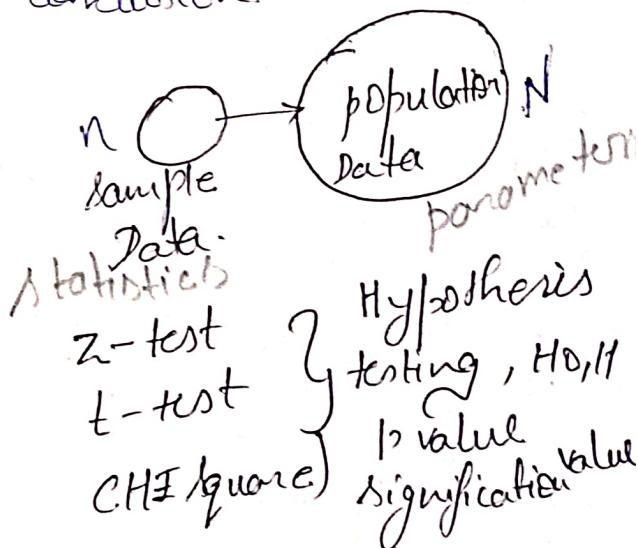
\* Measure of Central Tendency  
[mean, median, mode]

\* Measure of Dispersion  
[variance, std]

\* Different type of Distribution of data.  
[Histogram, pdf, pmf]

### ② Inferential Statistics

→ It consists of using data you have measured to form conclusion.

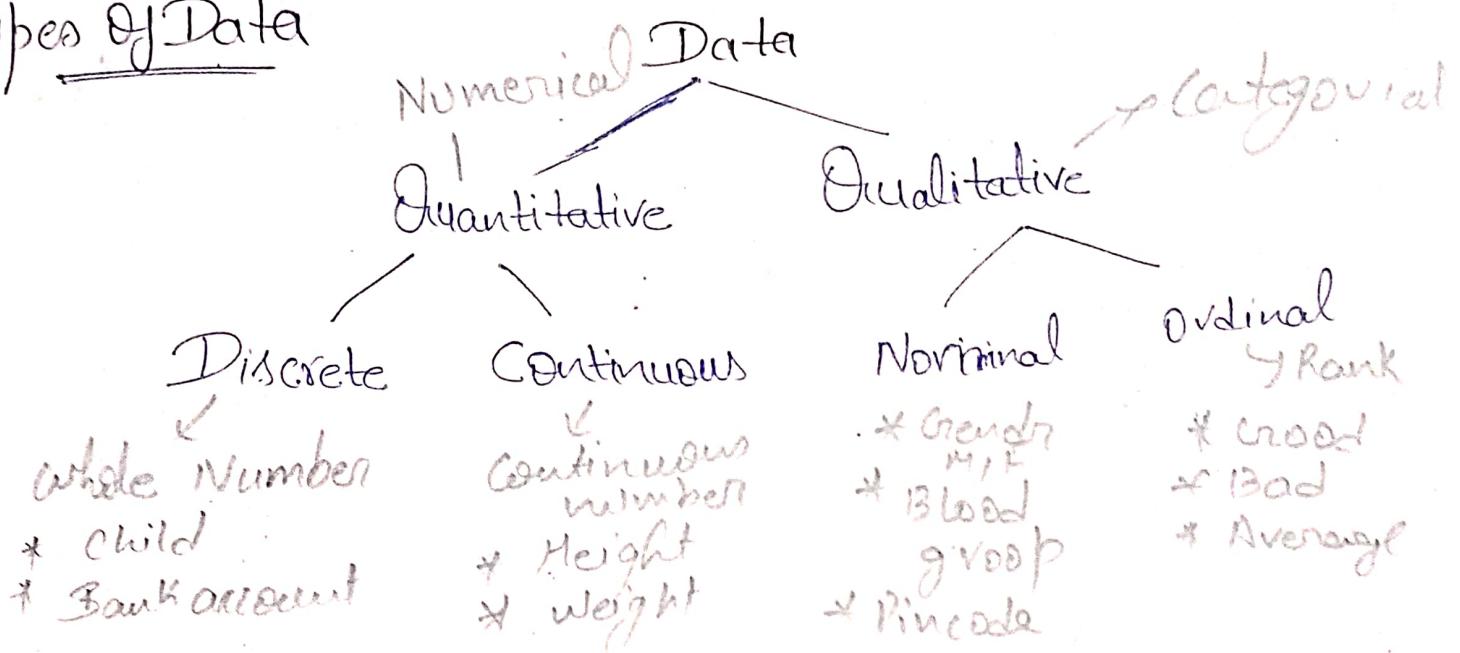


## Population and Sample Data.

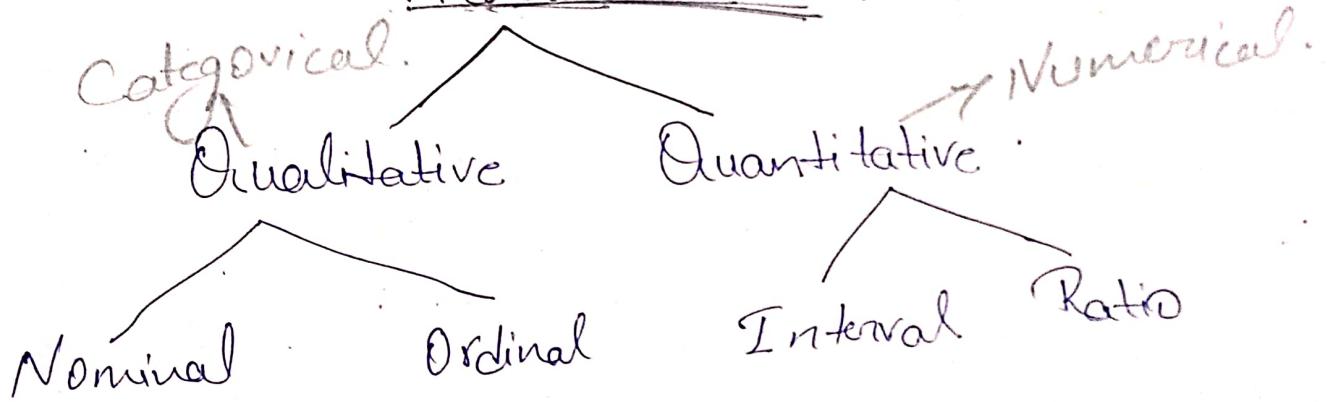
Population: the group you are interested in studying

sample: a subset of population.

# Types of Data



## Measurement level



### \* Nominal Scale Data

- \* Order does not matter
- \* Gender, colors, labels

eg.: favorite color  
 Red → 5 50%  
 Blue → 3 30%  
 Orange → 2 20%

### \* Ordinal Scale Data

- \* Ranking is important
- \* Order matters
- \* Difference cannot be measured

eg.: 1 → Best ↗ Diff  
 2 → Crood ↗ Diff  
 3 → Bad ↗ Diff

### \* Interval Data (context)

- \* The order matters
- \* Difference can be measured
- \* Ratio cannot be measured
- \* No true "0" starting point

OK → -2334J

eg.: Temperature

30°F ↗ 30°F  
 60°F ↗ 60°F

[3:1] X

120°F ↗ 120°F

- \* Ratio Scale Data (context) eg: student marks in class
  - \* The order matter  $0^{\circ}K \rightarrow 237^{\circ}C$
  - \* Difference are measurable (including ratio)
  - \* Contains a "0" starting point  $\rightarrow 0$  mark, starting no negative

## Measure of Central Tendency: (Descriptive statistics)

- \* Mean, \* Median, \* Mode

### (i) Mean:

Population ( $N$ )

$$\text{Population mean } (\mu) = \sum_{i=1}^n \frac{x_i}{N}$$

Sample ( $n$ )

$$\text{Sample mean} = \sum_{i=1}^n \frac{x_i}{n}$$

### (ii) Median:

\* Sort the Random variable.

\* No. of element if count is even then  $(\text{mid} + \text{mid} + 1)/2$   
else mid

- \* It is used to find the Central Tendency when outliers is present.

### (iii) Mode:

- \* frequency maximum
- \* for categorical  $[1, 1, 1, 2, 3, 3, 4, 5, 5] \rightarrow 1$

## 2) Measure of Dispersion [spread of the Data]

\* Variance \* standard deviation

### i) Variance

#### Population Variance

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

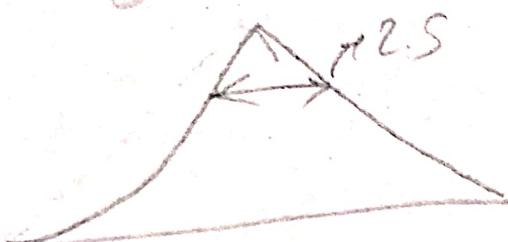
$x_i$  = Data points

$\mu$  = population mean

$N$  = Population size

Let!

$$\sigma^2 = 2.5$$



Variance                      Sample Variance.

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

$x_i$  = Data Points

$\bar{x}$  = Sample mean

$n$  = Sample size.

$(n-1)$  → Unbiased estimator  
to create.

$$s^2 = 7.5$$



### 2) Standard Deviation

#### Population Standard Deviation

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

$$\text{Let! } d \{1, 2, 3, 4, 5\}$$

$$\bar{x} = 3$$

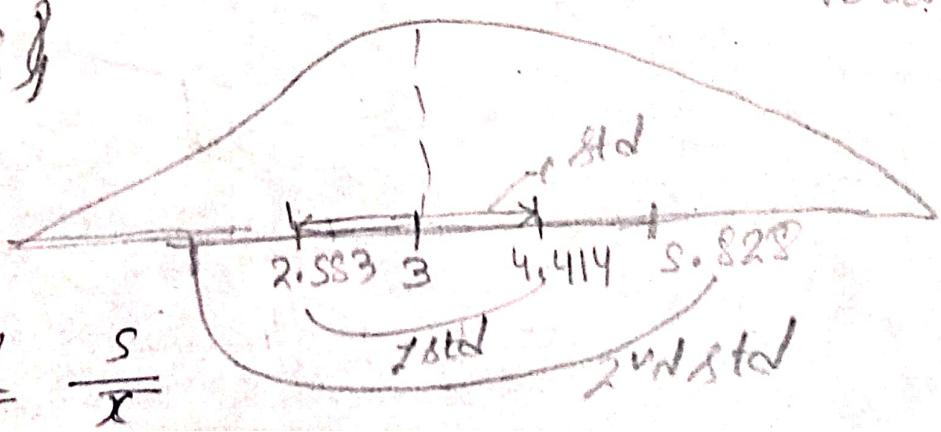
$$\sigma = 2.414$$

Sample Standard deviation  
Coefficient of variation =

#### Sample Std

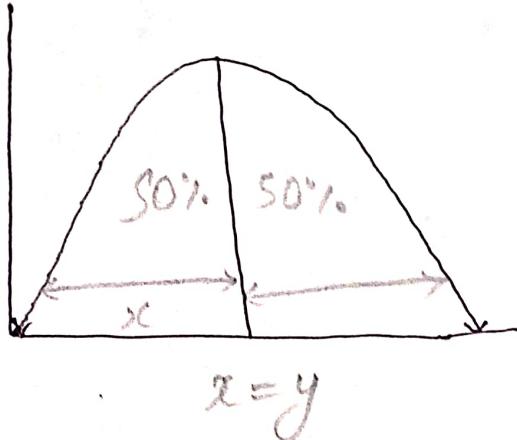
$$\text{std} = \sqrt{s^2}$$

$s^2$  = Sample Variance

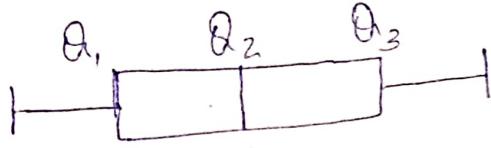


Skewness: measure of the asymmetry

## \* Normal / Gaussian / symmetrical Distribution

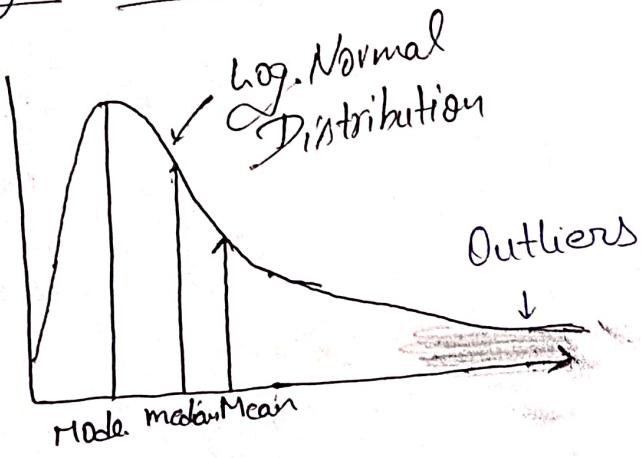


- \* No skewness
- \* Mean = Median = Mode

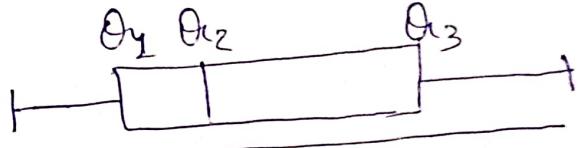


$$Q_3 - Q_2 \equiv Q_2 - Q_1$$

## \* Right Skewed



- \* Positive skewed
- \*  $\text{mean} \geq \text{median} \geq \text{mode}$
- \* outliers on right side

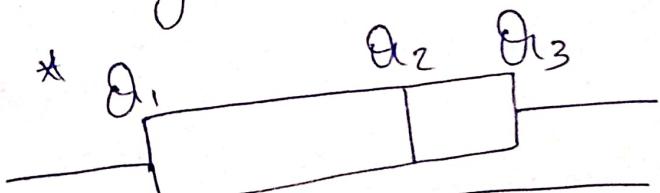


$$Q_3 - Q_2 \gg Q_2 - Q_1$$

## \* Left Skewed Distribution



- \*  $\text{mean} \leq \text{median} \leq \text{mode}$
- \* outliers on left side
- \* Negative skewed



$$Q_2 - Q_1 \gg Q_3 - Q_2$$

## Covariance (Direction only)

Covariance measure the direction of the linear relationship between two variable.

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

$$S_{xy} = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1} \rightarrow \text{Sample formula}$$

- \* Covariance gives a sense of direction
  - greater than  $> 0$  the two variables move together
  - less than  $< 0$  the two variables move in opposite directions
  - $= 0$  the two variables are independent.
- \* Covariance is not standardized, so its value depends on the scale of the variable e.g. ( $\rightarrow \infty$ )

## Correlation: *coefficient* Indicates direction & strength.

- \* Standardized (unit-free)
- \*  $-1 \rightarrow +1$
- \* Indicates direction & strength
- \* Use to compare relationships across datasets with different scales.

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

correlation range [-1 to 1]	
* 1	perfect positive corr
* -1	perfect negative
* 0	No linear relation

- Types:
- \* Pearson's  $\rightarrow$  same
  - \* Spearman Rank  $\rightarrow$   $\frac{\text{cov}(x,y)}{(\text{R}(x)) \cdot (\text{R}(y))}$

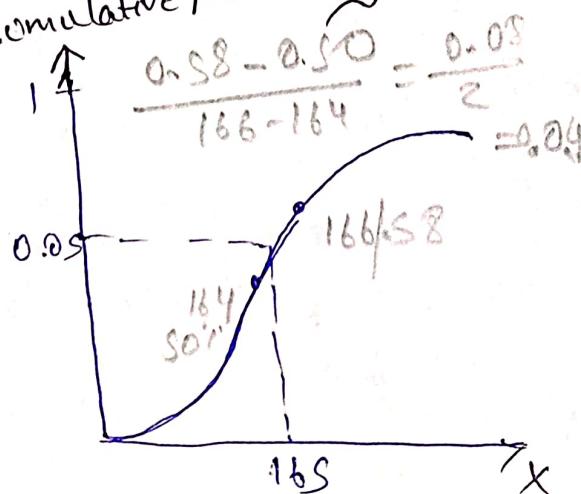
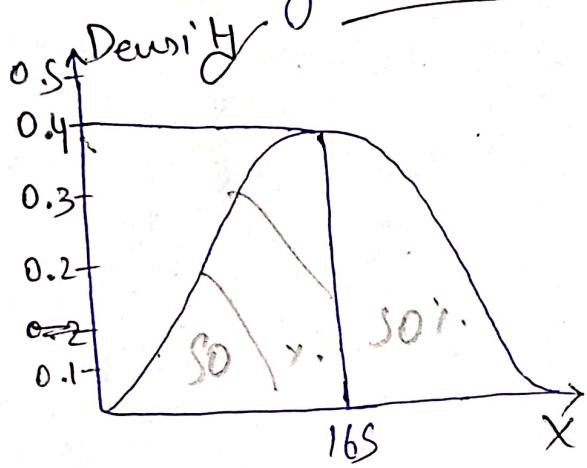
$\sigma_x, \sigma_y$  are the standard deviation of  $x$  and  $y$

# Probability Distribution function / Density function

Pdf / Pmf / CDF (Cumulative Distribution function)

\* Probability Density function (Pdf) / (CDF)

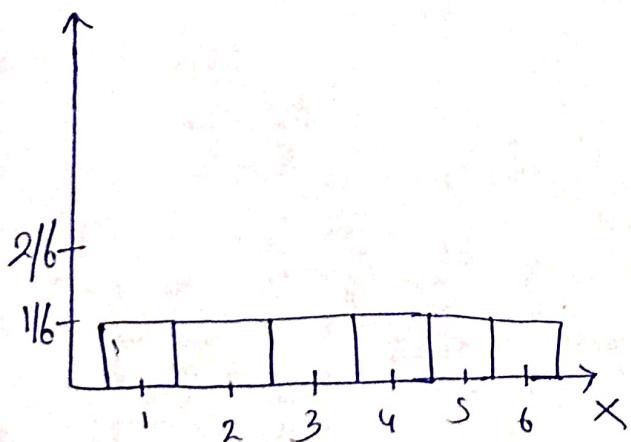
\* Distribution of Continuous Random Variable



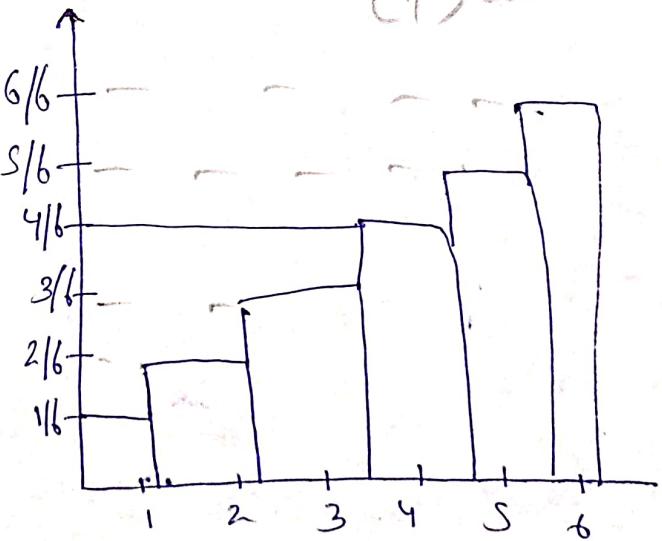
\* Probability Mass function:

\* Discrete Random Variable.

Eg: Rolling Dice



Cumulative probability  
(+) odd



# Types of Probability Distribution:

- (i) Normal / Gaussian Distribution (pdf) → Continuous
- (ii) Bernoulli Distribution (pmf) → Discrete / binary
- (iii) Uniform Distribution (pmf)
- (iv) Poisson Distribution (pmf)
- (v) Log Normal Distribution (pdf)
- (vi) Binomial Distribution (pmf)

## \* Bernoulli Distribution

\* Discrete Random Variable (pmf)

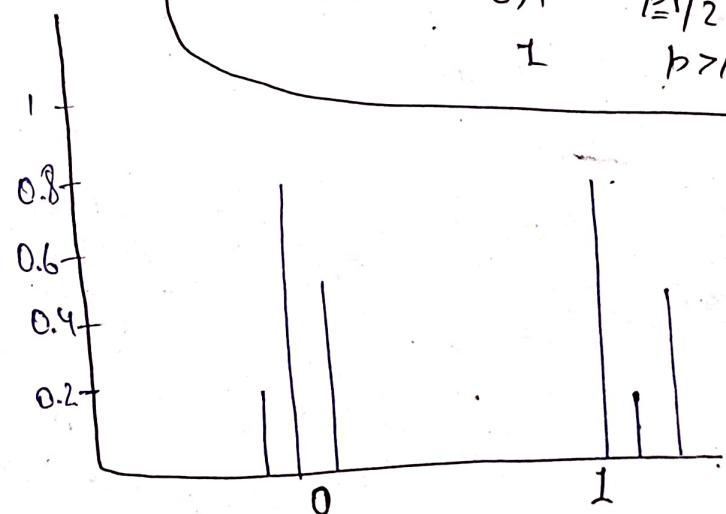
\* Outcomes are Binary {0, 1}

e.g.: Tossing a coin (H, T)

$$Pr(H) = 0.5 = p$$

$$Pr(T) = 0.5 = 1-p = q$$

\* Variance =  $pq$   
 \* Std =  $\sqrt{pq}$   
 \* Mean =  $p$   
 \* Median =  $\begin{cases} 0 & \text{if } p < 1/2 \\ 0,1 & \text{if } p = 1/2 \\ 1 & \text{if } p > 1/2 \end{cases}$



Let's whether the person will pass/fail

$$Pr(\text{Pass}) = 0.7 = p$$

$$Pr(\text{fail}) = 1 - 0.7 = 0.3$$

$$0.3 = q$$

$$1 - p = q$$

Three examples of Bernoulli Distribution

$$P(X=0) = 0.2 \text{ and } P(X=1) = 0.8$$

$$P(X=0) = 0.8 \text{ and } P(X=1) = 0.2$$

$$P(X=0) = 0.5 \text{ and } P(X=1) = 0.5$$

$$\text{PMF} = p^K \times (1-p)^{1-K} \quad \forall K = 0, 1$$

## \* Binomial Distribution: (pmf), $\text{EB}(n, p)$ , $n=1, \dots, n$

- \* A single success/failure experiment is also called a Bernoulli trial distribution and sequence of outcomes is said to Binomial.
  - \* Sequence of independent Bernoulli trials can collectively follow Binomial distribution.
  - \* Discrete Random Variable
  - \* Every experiment outcome is binary.
  - \* These experiments/trials are performed for  $n$  times
- eg!: Tossing a coin 10 times.
- 

Support:  $K \in \{0, 1, 2, 3, \dots, n\} \rightarrow$  Number of successes

PMF:  $\Pr(K, n, p) = {}^n C_K p^K (1-p)^{n-K}$

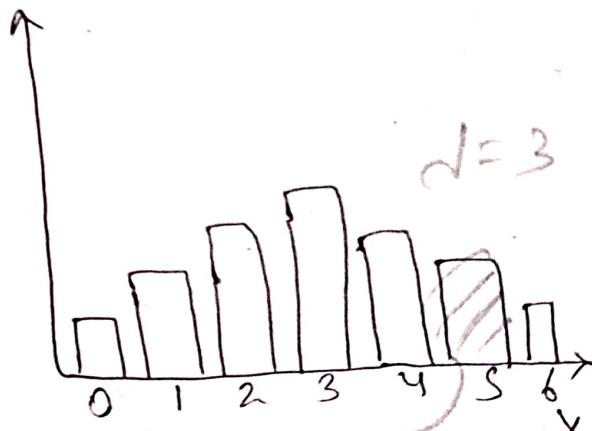
Mean $\rightarrow np$
Variance $\rightarrow npq$
Std $\rightarrow \sqrt{npq}$

## \* Poisson Distribution: (pmf)

- \* Discrete Random Variable

\* Describe the number of events occurring in a fixed time interval.

eg!: No. of people visiting hospital every hour.



\*  $\lambda \rightarrow$  Expected No. of events to occur at every time interval.

\*  $T =$  Time interval

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

$$P(X=5) = \frac{e^{-3} 3^5}{5!}$$

$$= 0.101 \\ = 10.1\%$$

$\lambda = 3$

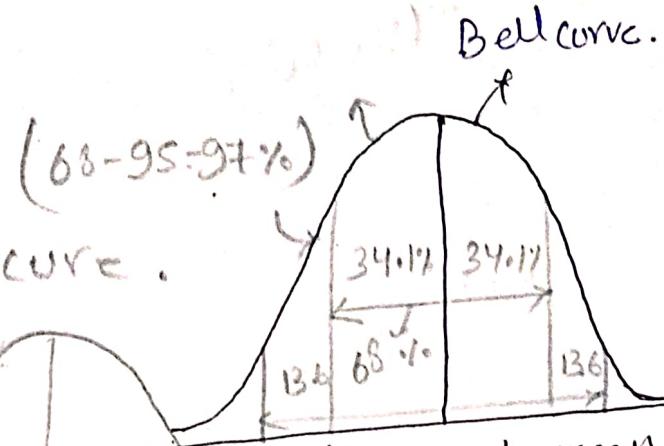
mean =  $\lambda \times t$

Variance =  $\lambda \times t$   
=  $\lambda \times t$

## \* Normal | Gaussian distribution (pdf),

### \* Continuous Random Variable

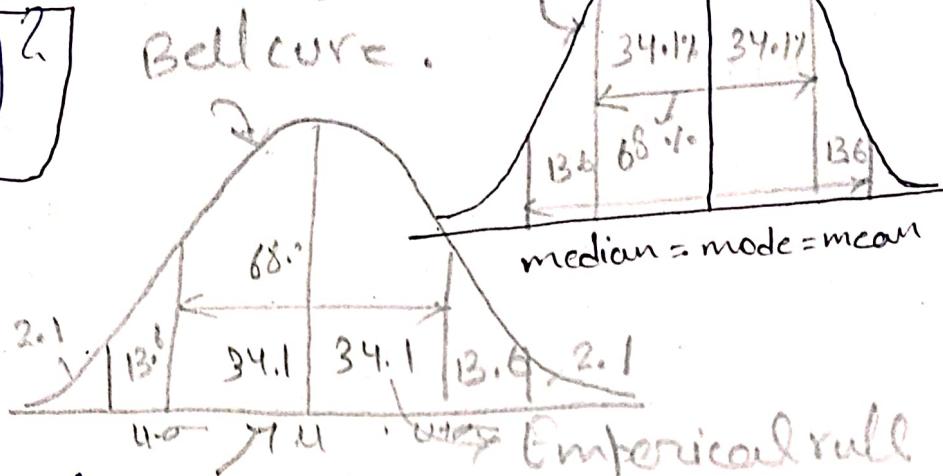
$$\text{P.D.F.} = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



\* mean =  $\mu$  = Average

\* Variance =  $\sigma^2$

\* Std =  $\sigma$



Empirical rule of Normal Distribution

$\Rightarrow (68-95-99.7)\%$ .

e.g.: Weight of the student in the class  
Height "

## \* Uniform Distribution

- \* Continuous Uniform Distribution (pdf)  $\rightarrow$  Random
- \* Discrete Uniform Distribution (pmf)

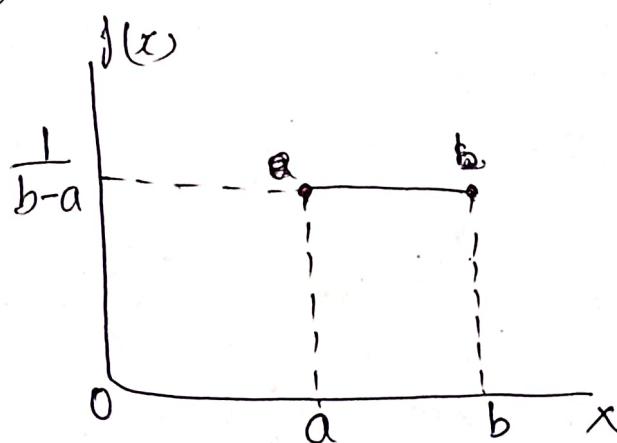
### Continuous Uniform Distribution (pdf)

\* Continuous random Variable

\* minimum & maximum values ( $a, b$ )

\*  $a < b$

$$\text{pdf} = \begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$



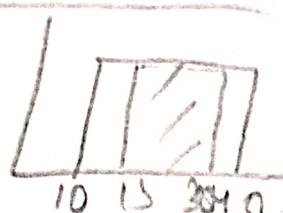
$$\text{Mean} = \frac{1}{2}(a+b)$$

$$\text{Median} = \frac{1}{2}(a+b)$$

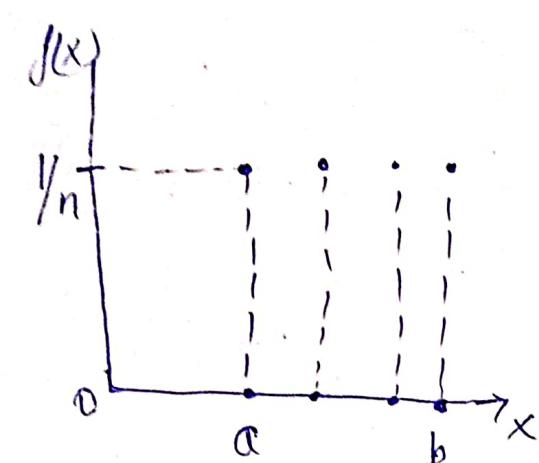
$$\text{Variance} = \frac{1}{12}(b-a)^2$$

(i) if max 40 & min 10  
(ii) fall between 15 & 30

$$P(15 \leq x \leq 30) = \frac{(x_2 - x_1)}{b-a} = 15 \times \frac{1}{30} = 50\%$$



## \* Discrete Uniform Distribution



\* finite number of outcomes equally likely to happen.

\* Discrete Uniform Distribution

e.g.: Rolling a dice  $\rightarrow \{1, 2, 3, 4, 5, 6\}$

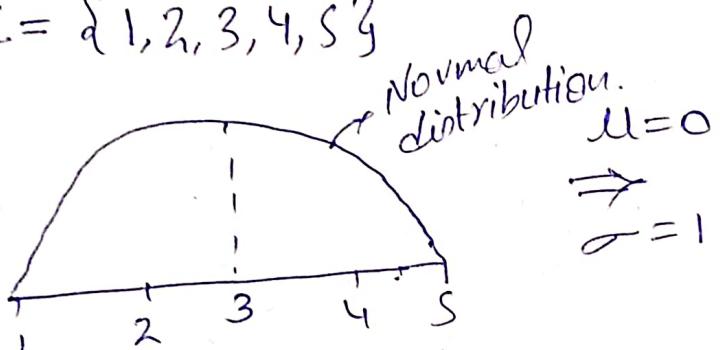
Notation  $\sim U(a, b)$

PMF  $\Rightarrow f(x) = \frac{1}{n}$

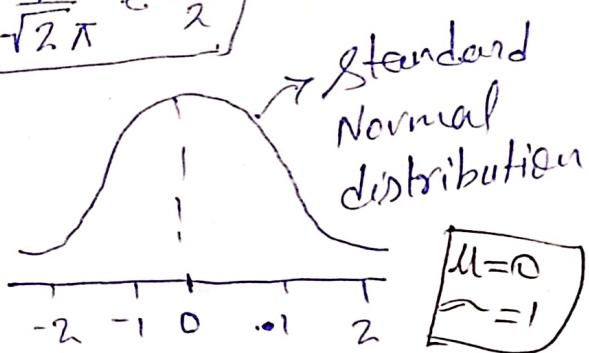
$$\text{mean, median} = \frac{a+b}{2}$$

## Standard Normal Distribution and Z-score

$$x = \{1, 2, 3, 4, 5\}$$

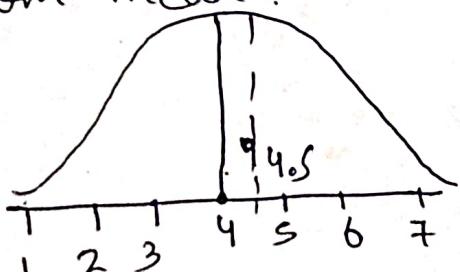


$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$



\* To convert Normal distribution to Standard Normal distribution we use Z-score.  $= \frac{x_i - \mu}{\sigma}$

Q1: How many student do stand and deviation 4.5 is always from mean?



$$x_1 = 4.5 \quad \sigma_0 = \frac{4.5 - 4}{1} = 0.5$$
$$\mu = 4 \quad \sigma = 1$$

Using Z-score

Q2: what percentage of data is falling above 4.5



$$\text{So, } z\text{-score} = \frac{4.5 - 4}{1} = 0.5$$

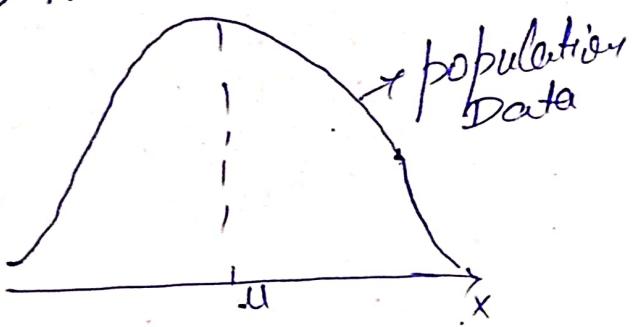
using Z-score table

$$= 1 - 0.69146$$
$$= 0.3085$$

$$= 30\%$$

Central Limit Theorem (CLT) + Always. Normal distribution

$$\textcircled{1} \quad X \sim N(\mu, \sigma)$$



falling  
sample (n<sub>7</sub>)

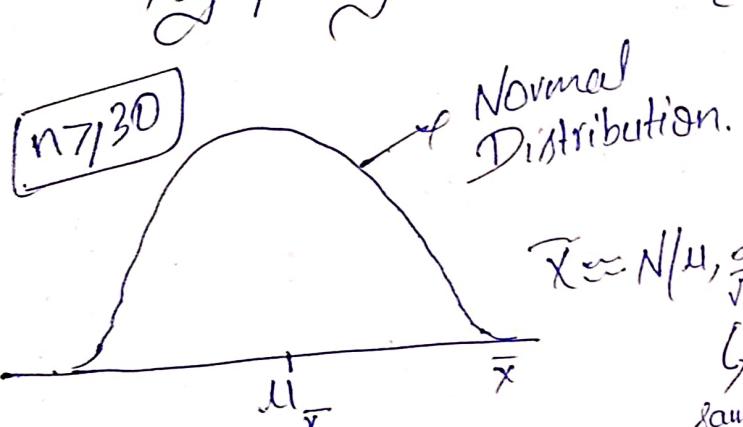
$$S_1 = \{ x_1, x_2, x_3, \dots \} = \bar{x}_1$$

$$S_2 = \{x_1, x_3, \dots\} = \bar{x}_2$$

83

$$S_{30} = \{ \_ \_ \_ \_ \_ \} = \overline{m}$$

So, By plotting that mean we get Normal distribution  
Note (even, poisson, binomial)



## Note

Note (even, possessive, no - 3)  
8. 1. 1. 1.

\* Population if Normal distribute

then N.o. sample can be any

If population is not nominal distributed then N.of. sample  $n > 30$

sample error of  
the mean

Estimate:

Estimate: It is an observed numerical value used to estimate an unknown population parameter.

### i) Point Estimate

① Single numerical value used to estimate the unknown parameters

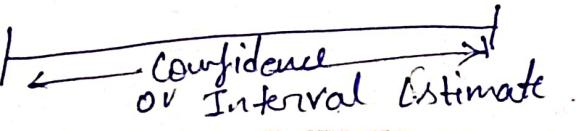
- \* Sample mean is a point estimate of a population mean.

$$C.I = \bar{x} + Z_{\alpha/2} \times \frac{\sigma}{\sqrt{n}} \quad \begin{matrix} \text{critical value} \\ \text{sample mean} \end{matrix} \quad \begin{matrix} \rightarrow \text{population} \\ \text{sample size.} \end{matrix}$$

④ Interval Estimate.

Range of values used to estimate the unknown population parameter

\* Interval estimate of population parameters are called Confidence Intervals



## P-value

The p-value is a number, calculated from a statistical test, that describes, how likely you are to have found a particular set of observations if the null hypothesis were true. P-values are used in hypothesis testing to help decide whether to reject the null hypothesis.

- \* Low p-value ( $< 0.05$ ): strong evidence against the null hypothesis  
→ Reject the null hypothesis
- \* High p-value ( $> 0.05$ ): weak evidence ag...  
→ fail to reject the null hypothesis

## Significance value!

Significance value is also known as alpha ( $\alpha$ ) ( $0.05, 0.01, 0.10$ )  
S.V., 1%, 10%

### Relationship b/w p-value & S.V.

- \* If  $p\text{-value} \leq \alpha$ , reject the null hypothesis
- \* If  $p\text{-value} > \alpha$ , fail to reject null hypothesis

## Confidence level!

The probability that the Confidence Interval contains the true population parameter.

$$C.I. = (1 - \alpha)$$

if: Significance value ( $\alpha$ ) = 0.05 then C.I. =  $1 - 0.05$   
 $= \underline{\underline{0.95}}$

## Error

Type 1 Error: We reject the Null Hypothesis when in reality it is true

Type 2 Error: We retain the null hypothesis when in reality it is false.

# Hypothesis Testing and Statistical Analysis

① Z-test  $\rightarrow$  Average

② T-test  $\rightarrow$

③ Chi Square  $\rightarrow$  Categorical data

④ ANNOVA  $\rightarrow$  Variance

Z-test is a statistical test used to determine whether there is a significant difference between sample and population mean or between two sample mean.

## Z-Test

\* population std  $\times n \geq 30$   
\* data normally distributed \* Sample should be independent  
\* population mean

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

Sample mean

Q1

The average heights of all residents in a city is 168cm. with  $\sigma = 3.9$ . A doctor believes the mean to be different. He measured the height of 36 individuals and found the average heights to be 169.5cm.

(a) State null and alternate hypothesis

(b) At a 95% confidence level, is there enough evidence to reject the null hypothesis?

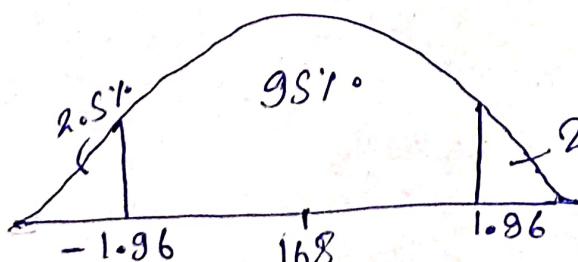
Sol:  $\mu = 168 \quad \sigma = 3.9, \bar{x} = 169.5, n = 36$

a) Null Hypothesis  $H_0: \mu = 168\text{cm}$

b) Alternate hypothesis  $H_1: \mu \neq 168\text{cm}$  & 2 Tail Test

c) C.I = 0.95  $\alpha = 1 - 0.95 = 0.05$

Note! & things  
\* Z-test  
\* p-value



area under the curve =  $1 - 0.05 = 0.95$

Using Z-score table we find (0.95)

$$\text{Z-score} = 1.96$$

d) Statistical Analysis / formulae apply

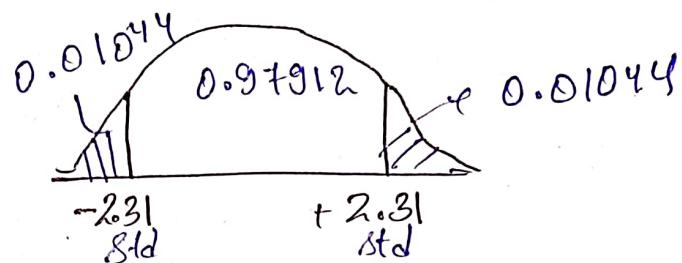
$$Z\text{-score} = \frac{x_i - \mu}{\sigma/\sqrt{n}} \text{ for population}$$

$$Z\text{-test} = \frac{x_i - \mu}{\sigma/\sqrt{n}} \Rightarrow \frac{169.5 - 168}{3.9/\sqrt{36}} \Rightarrow 2.31$$

If Z-test value is less than -1.96 or greater +1.96  
then we Reject the Null Hypothesis  
else, we Accept Null Hypothesis

$2.31 > +1.96$  & we Reject the Null Hypothesis

② P-value (By using)



$$Z\text{-test} = \frac{169.5 - 168}{3.9/\sqrt{36}} \Rightarrow 2.31$$

$$= 1 - \text{area under the curve of } +2.31 \Rightarrow 1 - 0.97912 \Rightarrow 0.01044$$

$$p\text{-value} = 0.01044 + 0.01044 = \underline{\underline{0.02088}}$$

if p-value < significance

$0.02088 < 0.05$  & we Reject the Null Hypothesis

② A factory manufactures bulbs with a average warranty of 5 years with standard deviation of 0.50. A worker believes that the bulb will malfunction in less than 5 years. He tests a sample of 40 bulbs and find the average time to be 4.8 years.

a) State null and alternate hypothesis

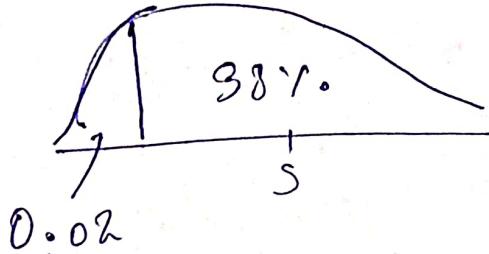
b) At a 2% significance level, is there enough evidence to support idea that the warranty should be revised?

Sol.  $\mu = 5$ ,  $\sigma = 0.50$ ,  $n = 40$ ,  $\bar{x} = 4.8$

1) Null Hypothesis  $\Rightarrow H_0 : \mu = 5$

2) Alternative Hypothesis  $H_1 : \mu < 5$  & 1 tail Test  
 C.I = 0.98      or  
 $C.I = 1 - 0.02 = 0.98$   
 $\alpha = 1 - 0.98 = 0.02$

3) Decision Boundary



Using z-score of -2.53 we get. 0.0570



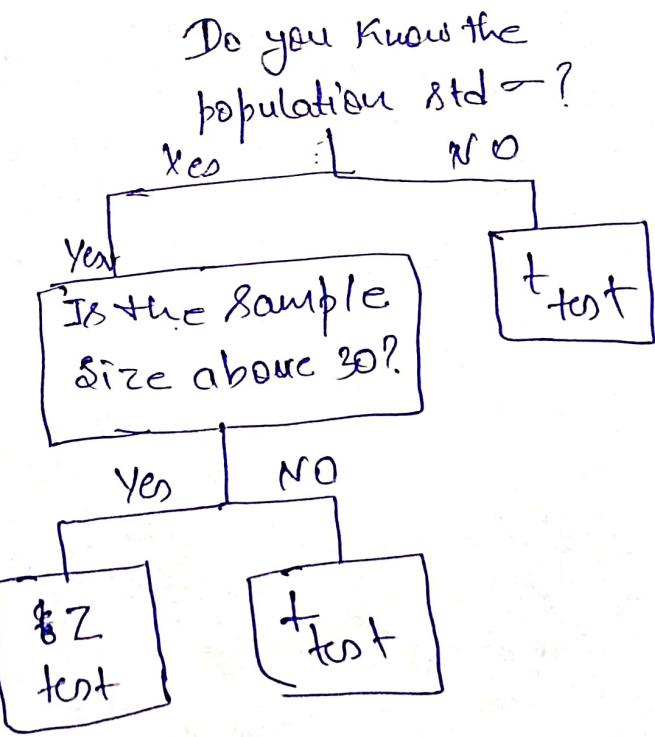
4) P-value  $\rightarrow$  Z-test =  $\frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$   
 $= \frac{4.8 - 5}{0.50 / \sqrt{40}}$   
 $= -2.53$

Area under the curve of -2.53, Z-value is = 0.0570

p-value  $\rightarrow$  0.0570

$\therefore$  if p-value is significant  
 $0.0570 < 0.02 \Rightarrow$  false  $\therefore$  We accept the Null Hypothesis

T-test	Z-test
Small ( $n \leq 30$ )	Large ( $n > 30$ )
Population variance unknown	Know
Distribution skewed Student's t-distribution	Normal distribution



## T-test

A t-test is a statistical test used to compare the mean of one or more groups when the population standard deviation is unknown and the sample size is small ( $n < 30$ ).

### \* Assumption:

- \* The data should be numerical, not categorical.
- \* The sample must be randomly selected: Random Sampling
- \* The data should follow an approximately normal distribution.

### formula

#### (a) One-Sample

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

$s$  = sample std

$\mu$  = population mean

$\bar{x}$  = sample mean

$n$  = population size

#### (b) Two Sample

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$\bar{x}_1, \bar{x}_2$  = sample mean

$s_1^2, s_2^2$  = sample variance

$n_1, n_2$  = sample size

#### (c) Paired T-Test

$$t = \frac{\bar{d}}{s_d/\sqrt{n}}$$

$\bar{d}$  = Mean of difference  
b/w paired values

$s_d$  = std of difference

$n$  = Number of pairs

### Degrees of freedom (df)

The Degrees of freedom (df) represent the number of independent values in a statistical calculation that are free to vary.

#### (a) One Sample

$$df = n - 1$$

$n$  = sample size

#### (b) Two Sample

$$df = n_1 + n_2 - 2$$

$n_1, n_2$  = sample sizes  
of two groups

#### (c) Paired T-Test

$$df = n - 1$$

$n$  = Number of paired observations.

\* Df determine the shape of t-distribution

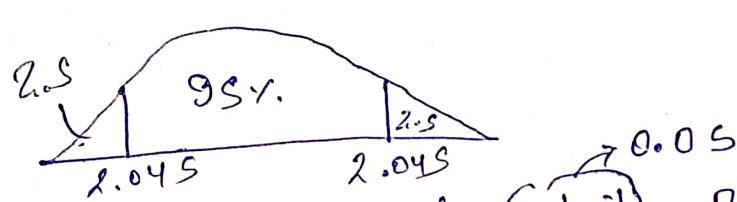
\* Smaller d.f. mean wider distribution. \* larger closer to normal

## A. t. Test

In the population the average IQ is 100. A team of researcher want to test a new medication to see if it has either a positive or negative effect on intelligence or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140 with a standard deviation of 20. Did the medication affect intelligence? C.I = 95%

$$\text{Sol: } \mu = 100, n = 30, s = 20, \bar{x} = 140, CI = 95\% \\ \alpha = 0.05$$

- ① Null Hypothesis  $H_0: \mu = 100$
- ② Alternate Hypothesis  $H_1: \mu \neq 100$  of 2 tail Test
- ③  $\alpha = 0.05$
- ④ Degree of freedom =  $n - 1 = 30 - 1 = 29$
- ⑤ Decision Rule



By using t-table and the d.f is 29 for 2 tail = 2.045  
t.table 29 and 2 tail = 0.05  
= 2.045

If t test is less than -2.045 and greater than 2.045, then Reject the Null Hypothesis

- ⑥ calculate t test stat

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{140 - 100}{20/\sqrt{30}} = 10.96$$

- ⑦ conclusion We reject the null hypothesis

$$2.04 < 10.96$$

# Confidence Interval and Margin of Error



We construct a confidence interval to help estimate what the actual value of the unknown population mean is.

$$\text{Point Estimate} \pm \text{Margin of Error}$$

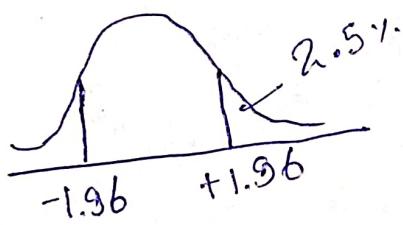
$$Z \text{ test} = \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

Q. In the verbal section of CAT exam the standard deviation is known to be 100. A sample of 25 test taken has a mean of 520. Construct a 95% C.I about the mean.

Sol  $\bar{x} = 520, \sigma = 100, n = 25, C.I = 0.95, \alpha = 0.05$

$$\alpha = 0.05$$

$$\begin{aligned} &= 1 - 0.05 \\ &= 0.975 \rightarrow \text{(using z-table)} \\ &= 1.96 \end{aligned}$$



$$\bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$\text{Lower C.I} = 520 - (1.96) * \frac{100}{\sqrt{25}} = 480.8$$

$$\text{Higher C.I} = 520 + (1.96) * \frac{100}{\sqrt{25}} = 559.2$$



So range between 480.8 to 559.2

So, I am 95% confident that the mean CAT score lies between 480.8 to 559.2

## Chi-Square Test.

The Chi-Square Test help us check if two things are related or if the difference we see are just luck.

e.g./ you ask 100 people if they like pizza or burgers, and you check if men and women have different preferences  
the chi-square Test you if gender really affects their choice or if the difference happened by chance.

## Types of Chi-Square Test:

- |  |  |
|--|--|
| (i) Chi square Goodness of-fit Test<br>* 1 categorical variable<br>* Tests if a sample follow a known distribution<br>* Use: you compared observed vs expected frequencies | (ii) Chi-Square Test of independence<br>* 2 categorical variable.<br>* Tests if two categorical variable are related or independent<br>* you check if one category affects the others. |
|--|--|

formula:  $\chi^2 = \sum \frac{(O - E)^2}{E}$

- Q. In 2010 Census of the City the weight of the individuals in a small city were found to be the following.
- |         |       |     |
|---------|-------|-----|
| ≤ 50 kg | 50-75 | 775 |
| 20%     | 30%   | 50  |

In 2020 ages of  $n=500$  individuals were sampled. Below are result.

≤ 50	50-75	775
140	160	200

Using  $\alpha=0.05$ , would you conclude the population difference of weights has changed in the last 10 years.

<50	50-75	>75
Expected		
20%	30%	50%

<50	50-75	>75
Observed		
140	160	200

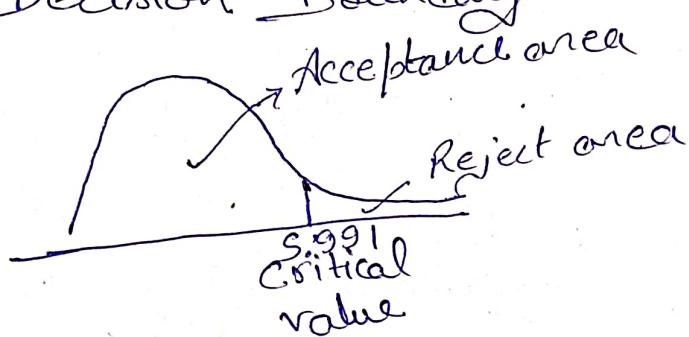
Expected	<50	50-75	>75
Expected	$\frac{20 \times 500}{100}$	$\frac{30 \times 500}{100}$	$\frac{50 \times 500}{100}$
	100	150	250

- ① Null Hypothesis!  $H_0$ : The data meets the exception  
 ② Alternative Hypothesis!  $H_1$ : The data does not meet the exception.

③  $\alpha = 0.05$  CI = 95%

④ Degree of freedom  $df = k-1 = 3-1 = 2$

⑤ Decision Boundary



By using Chi-table  
 $df = 2$  &  $\alpha = 0.05$   
 we get S.991

⑥ Calculate Chi square test statistics  
 $\chi^2 = \sum \frac{(O-E)^2}{E} \Rightarrow \frac{(40)^2}{100} + \frac{(10)^2}{150} + \frac{(50)^2}{250} = 26.66$

⑦ Conclusion  
 If  $\chi^2$  is greater than S.99; Reject  $H_0$   
 else we fail to reject the Null Hypothesis.

So,  
 $\chi^2 = 26.66 > 5.99$

{ Reject the Null Hypothesis? }

## F-Distribution

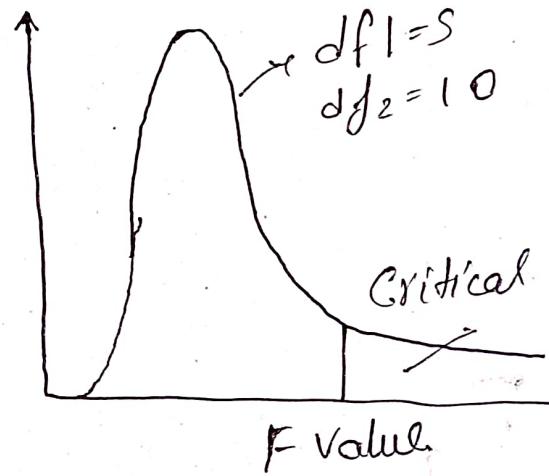
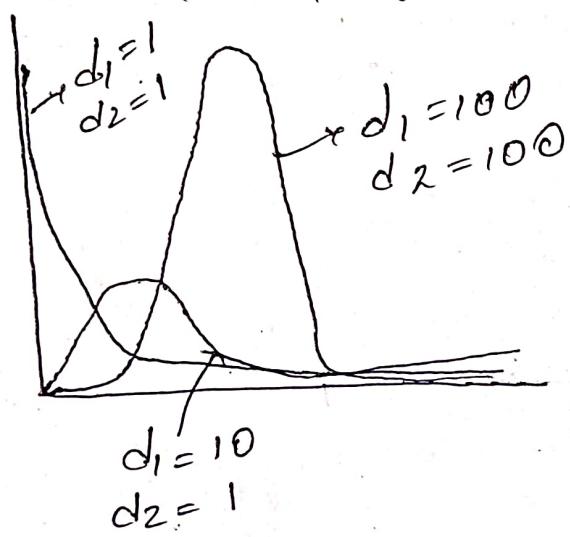
The F-Distribution is a probability distribution used mainly in ANOVA (Analysis of Variance) and F-Test to compare two datasets. It helps determine if the difference in data are real or just due to chance.

mean or variance  
by ANOVA

### Key features of F-Distribution

- \* Always positive (values are never negative)
- \* Right-skewed (not symmetric like a normal distribution)
- \* Used in hypothesis testing when comparing two variance
- \* Defined by two degrees of freedom ( $d_f 1, d_f 2$ )  $> 0$  always

### Probability density function



### formula

$$P_{df} = f(x; d_2, d_1) = \frac{\sqrt{(d_1 x)^{d_1} d_2^{d_2}}}{(d_1 x + d_2)^{d_1 + d_2}} \cdot \frac{x}{B\left(\frac{d_1}{2}, \frac{d_2}{2}\right)} \rightarrow \text{Beta function or } B(m, n)$$

$$B(m, n) = \frac{(m-1)! (n-1)!}{(m+n-1)!} = \frac{m+n}{mn} / \frac{(m+n)}{m}$$

$$F = \frac{s_1^2}{s_2^2} \therefore s_1^2, s_2^2 \text{ Variance of first and second sample.}$$

## F-Test

The F-Test is a statistical test used to compare the variance of two datasets. It helps check if the spread (variability) of data in two groups is significantly different.

Q. The following data shows the no. of bulbs produced daily for some days by 2 workers A and B.

A	40	30	38	41	38	35
---	----	----	----	----	----	----

B	39	38	41	33	32	39	40	34
---	----	----	----	----	----	----	----	----

Can we consider based on the data workers B is more stable and efficient.  $\alpha = 0.05$  C.I = 95% = 0.95

- Sol: ① Null Hypothesis  $H_0: \sigma_1^2 = \sigma_2^2$   
 ② Alternate Hypothesis  $H_1: \sigma_1^2 \neq \sigma_2^2$

③ calculate of variance

$$\begin{aligned} A & \text{ e. } (x_i - \bar{x})^2 \\ 40 & 37 9 \\ 30 & 37 49 \\ 38 & 37 1 \\ 41 & 37 16 \\ 38 & 37 1 \\ 35 & 37 4 \\ \hline \bar{x}_1 & = 37 \end{aligned}$$

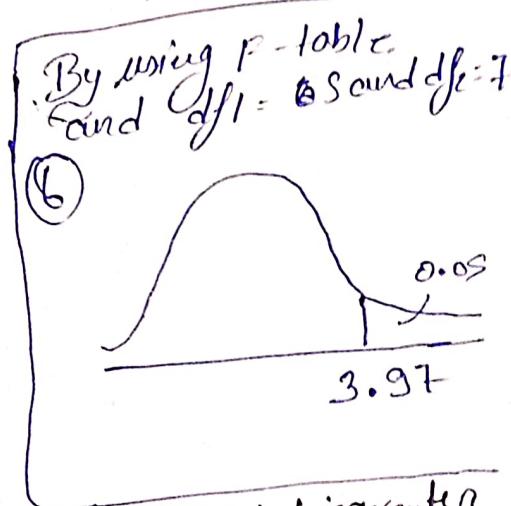
$$\begin{aligned} \sum (x_i - \bar{x})^2 &= 80 \end{aligned}$$

$$\begin{aligned} s_1^2 &= \frac{n}{n-1} \sum (x_i - \bar{x}_1)^2 \\ &= \frac{80}{6-1} \\ &= 13 \end{aligned}$$

$$\begin{aligned} B & \bar{x}_2 \\ x_2 & 39 37 38 41 37 33 32 37 39 40 37 34 \\ & \bar{x}_2 = 37 \end{aligned}$$

$$\begin{aligned} (x_2 - \bar{x}_2)^2 & 4 1 16 36 1 16 25 4 9 \\ \sum (x_2 - \bar{x}_2)^2 & = 84 \end{aligned}$$

$$\begin{aligned} s_2^2 &= \frac{84}{8-1} \\ &= \frac{84}{7} \\ &= 12 \end{aligned}$$



So, if F-test is greater than 3.97 it rejected the Null Hypothesis  
 $1.33 < 3.97$   
 we fail to reject Null

④ Calculate of variance ratio of F-test

$$F = \frac{s_1^2}{s_2^2} = \frac{13}{12} = 1.33$$

⑤ Decision Rule

$$df_1 = 6-1 = 5 \text{ and } df_2 = 8-1 = 7$$

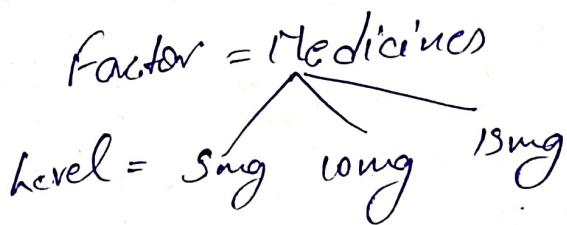
# ANOVA (Analysis of Variance)

ANOVA is a statistical method used to compare the mean of 2 or more groups

## Anova

- ① Factors (Variables)
- ② Level

eg:



## Assumption in ANOVA

- (i) Data is numerical and normally distributed
- (ii) Homogeneity of Variance  
Each one of the population has same variance  
 $\left[ \sigma_1^2 = \sigma_2^2 = \sigma_3^2 \right]$   
population variance in different levels of each independent variable are equal.
- (iii) Samples are independent and random
- (iv) Absence of outliers.

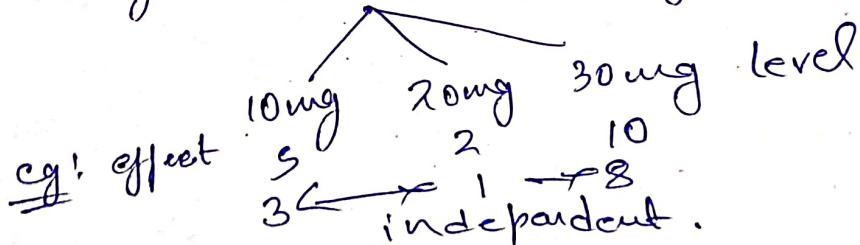
formula

$$F = \frac{\text{Variance between groups}}{\text{Variance within groups}}$$

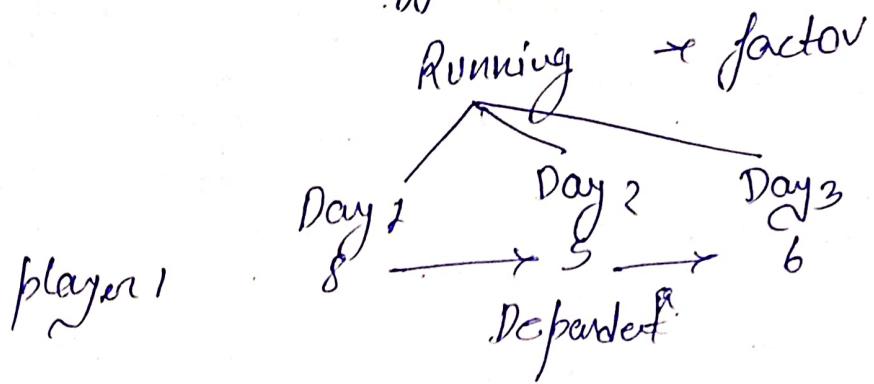
## Types of ANOVAs

- (i) One way ANOVA: 1 factor with atleast 2 levels those level are independent.

eg: Medication  $\rightarrow$  factor



- ② Repeated Measures ANOVA:  
 One factor with at least 2 levels, <sup>and level</sup> are dependents.
- \* It used when you measure the same subjects multiples times under different conditions or over time.



- ③ Factorial ANOVA / Two factor
- Two factor or more factor each of which with at least 2 level, level can be either independent and dependent

Running - Factor 1

	Day 1	Day 2	Day 3	
factor 2	Men	8	5	6
	Independent	7	4	3
	Women	6	5	4
	3	2	1	dependent.

### Hypothesis Testing in ANOVA

Null Hypothesis  $H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_R$

Alternate hypothesis  $H_1$ : At least one of the mean is not equal.

$$+ \boxed{\mu_1 \neq \mu_2 = \mu_3 = \dots = \mu_R}$$

Q.1 Doctors want to test a new medication with reduce headache. They splits the participant into 3 condition [15mg, 30mg, 45mg]. later on the doctor ask the patient to rate the headache between [1-10]. Are they any differences between the 3 condition using alpha 0.05

Sol: ① Null Hypothesis:  $H_0: \mu_{15} = \mu_{30} = \mu_{45}$

② Alternative Hypothesis  $H_1$ : not all  $\mu_i$  are equal.

$$\alpha = 0.05, C.I = 0.95$$

③ calculate Degree of freedom

$$N = 21 \quad a = 3 \quad n = 7$$

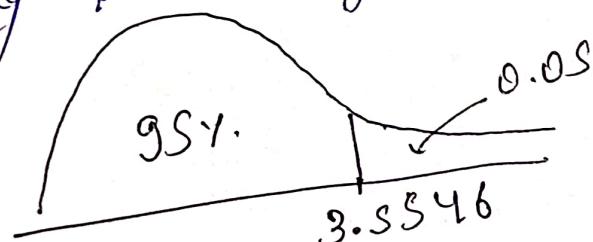
15mg	30mg	45mg
9	7	4
8	6	3
7	6	2
8	7	3
8	8	4
9	7	3
8	6	2

$$df_{\text{between}} = a - 1 = 3 - 1 = 2$$

$$df_{\text{within}} = N - a = 21 - 3 = 18$$

$$df_{\text{total}} = N - 1 = 21 - 1 = 20 \quad (\text{df}_{\text{between}} + \text{df}_{\text{within}})$$

Using F-table of (2, 18) we get = 3.5546  
Critical = 3.5546



④ State Decision Rule

If F is greater than 3.5546, reject the Null Hypothesis

## ⑤ Calculate test statistics

	SS	df	MS	F
Between	98.67	2	$\frac{98.67}{2} = 49.34$	$\frac{49.34}{0.54}$
within	10.29	18	$\frac{10.29}{18} = 0.54$	$= 86.56$
Total	108.96	20	49.88	

SS<sub>between</sub>, SS<sub>within</sub>, SS<sub>total</sub>

$$\text{SS}_{\text{between}} = \frac{\sum (\bar{x}_{ai})^2}{n} - \frac{T^2}{N}$$

$$15\text{mg} = 9+8+7+8+8+9+8 = 57$$

$$30\text{mg} = 7+6+6+7+8+7+6 = 47$$

$$45\text{mg} = 4+3+2+3+4+3+2 = 21$$

15mg	30mg	45mg
9	7	4
8	6	3
7	6	2
8	7	3
8	8	4
9	7	3
8	6	2

$$\text{SS}_{\text{between}} = \frac{57^2 + 47^2 + 21^2}{7} - \frac{[57 + 47 + 21]}{21} = 98.67$$

$$\text{SS}_{\text{within}} = \sum y^2 - \frac{\sum (\bar{x}_{ai})^2}{n}$$

$$= \sum y^2 - \left[ \frac{57^2 + 47^2 + 21^2}{7} \right]$$

$$\begin{aligned} \sum y^2 &= 9^2 + 8^2 + 7^2 + 8^2 + 8^2 + 9^2 + 8^2 + 7^2 + 6^2 + \dots \\ &= 853 \end{aligned}$$

$$= 853 - \left[ \frac{57^2 + 47^2 + 21^2}{7} \right]$$

$$\boxed{\text{SS}_{\text{within}} = 10.29}$$

$$\textcircled{S} \quad \underline{SS_{\text{total}}} = \sum y^2 - \frac{T^2}{N}$$

$$= 883 - \frac{123^2}{21} = 108.9S \quad (SS_{\text{between}} + SS_{\text{within}})$$

$$\textcircled{S} \quad F = \frac{MS_{\text{between}}}{MS_{\text{within}}}$$

$$= \frac{49.34}{0.84} = 86.56$$

So, con<sup>u</sup>  
 If  $F$  is greater than  $3.8846$ , reject the Null Hypothesis  
 $\boxed{86.56 > 3.8846}$  Reject the Null Hypothesis