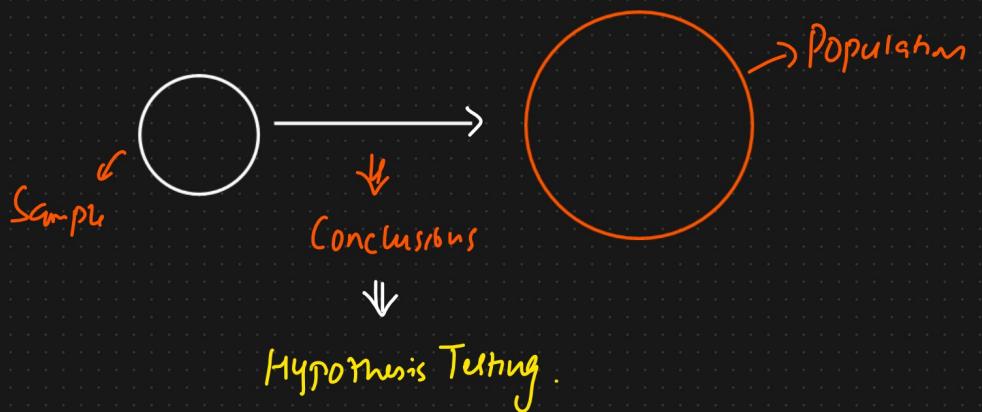


## A Hypothesis And Hypothesis Testing Mechanism

Inferential Stats  $\div$  Conclusion or Inference



Person  $\rightarrow$  Crime  $\rightarrow$  Court

### Hypothesis Testing Mechanism

- ① Null Hypothesis ( $H_0$ ) - Person is not guilty
  - The assumption you are beginning with
- ② Alternate Hypothesis ( $H_1$ ) - The person is guilty
  - Opposite of Null Hypothesis
- ③ Experiment  $\rightarrow$  Statistical Analysis  $\{P \text{ value, Significance value}\}$ .
  - $\rightarrow$  Direct Proof (DNA, Finger Test)
- ④ Accept the Null Hypothesis or Reject the Null Hypothesis

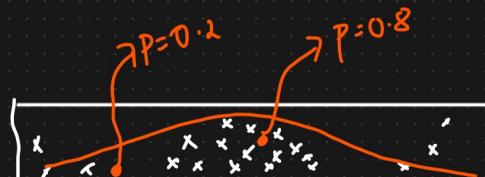
Eg: Colleges at District A states its <sup>Average</sup> passed percentage of Students are 85%. A new college opened in the district And it was found that a sample of student 100 have a pass percentage of 90% with a standard deviation of 4%. Does this college have a different passed percentage.

Ans) Null Hypothesis (H<sub>0</sub>) =  $\mu = 85\%$ .

Alternate Hypothesis (H<sub>1</sub>) =  $\mu \neq 85\%$ .

## 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.



Out of 100 touches, we touch around 20 times in this region

Hypothesis Testing Eg : Coin is Fair or Not {100 times}

$$P(H) = 0.5 \quad P(T) = 0.5$$

① Null Hypothesis :

$$H_0 : \text{Coin is fair} \quad P(H) = 0.6 \quad P(T) = 0.4$$

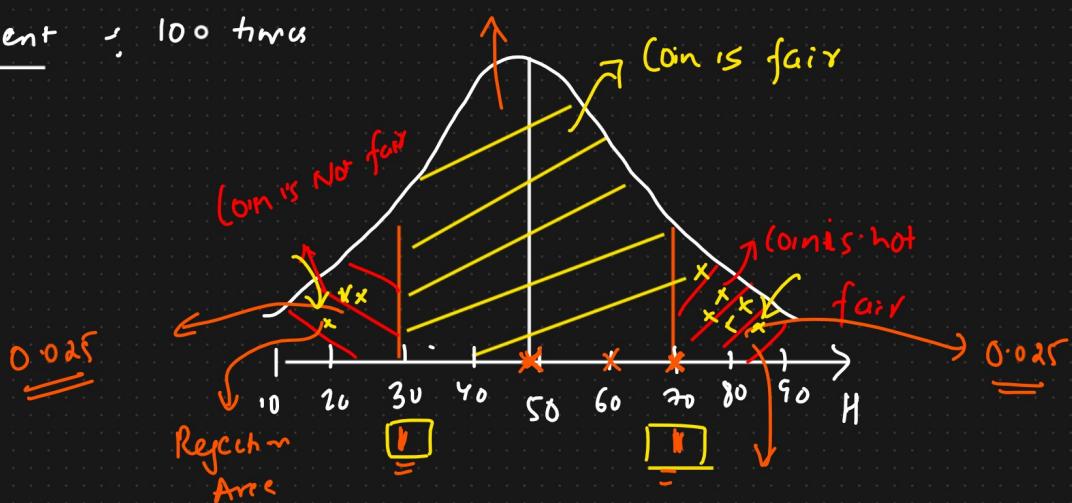
$$P(H) = 0.7 \quad P(T) = 0.3$$

② Alternative Hypothesis :

$$H_1 : \text{Coin is not fair}$$

95% C.I

③ Experiment : 100 times



④ Significance Value :  $\alpha = 0.05 \Leftarrow$

Rejection Area

$$C.I = 1 - 0.05 = 0.95$$

⑤ Conclusion  $P < \text{Significance value}$

Reject the Null Hypothesis

else

Fail to Reject the Null Hypothesis

## Hypothesis Testing And Statistical Analysis

- ① Z Test }  $\Rightarrow$  Average  $\Rightarrow$  Z table  $\rightarrow$  Z score And p value
- ② t Test  $\Rightarrow$  t table
- ③ CHI SQUARE  $\Rightarrow$  Categorical Data
- ④ ANNOVA  $\Rightarrow$  Variance

Z test. i) population std - ii)  $n \geq 30$

With a  $\sigma = 3.9$

i) The average heights of all residents in a city is 168cm. A doctor believes the mean to be different. He measured the height of 36 individuals and found the average height to be 169.5 cm.

(a) State null and Alternate Hypothesis

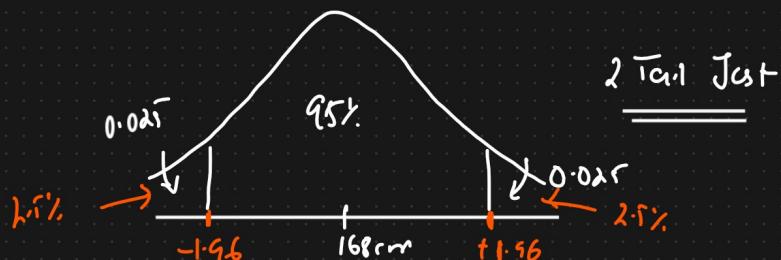
(b) At a 95% confidence level, is there enough evidence to reject the null hypothesis.

Ans)  $M = 168\text{cm}$   $\sigma = 3.9$   $n = 36$   $\bar{x} = 169.5$   
 $C.I = 0.95$   $\alpha = 1 - C.I = 1 - 0.95 = 0.05\%$

① Null Hypothesis  $H_0 = M = 168\text{cm}$

② Alternate Hypothesis  $H_1 = M \neq 168\text{cm}$

③ Based on C.I we will draw Decision Boundary



$$1 - 0.025 = 0.975 \Rightarrow Z\text{-score}$$

$$\Downarrow \\ \text{Area} \Rightarrow +1.96$$

if  $Z$  is less than  $-1.96$  or greater than  $+1.96$ , Reject the Null Hypothesis.

### Z-test

$$Z\text{-score} = \frac{\bar{X} - \mu}{\sigma}$$

$$\Downarrow \\ \frac{\sigma}{\sqrt{n}}$$

$$Z_d = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

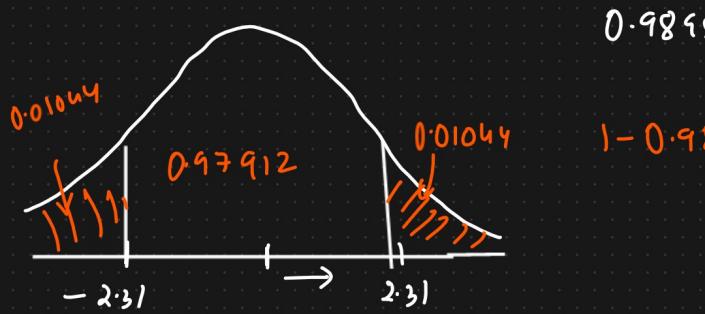
$$= \frac{169.5 - 168}{3.9/\sqrt{36}}$$

$$Z_d = \frac{1.5}{0.65} \approx 2.31$$

Conclusion

$Z\text{-score} \downarrow$   
 $2.31 > 1.96$  Reject the Null Hypothesis

$$P < 0.05$$



$$0.98956$$

$$1 - 0.98956 = 0$$

Final Conclusion the Average  $\neq 168\text{cm}$

The average height seems to increasing based on sample height.

$$\textcircled{1} \quad p\text{ value} = 0.01044 + 0.01044$$

$$= 0.02088$$

$$P < 0.05$$

$0.02088 < 0.05 \Rightarrow$  Reject the Null Hypothesis

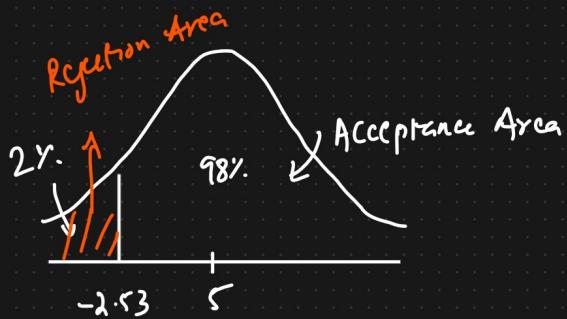
② A factory manufactures bulbs with a average warranty of 5 years with standard deviation of 0.50. A worker believes that the bulb will malfunction in less than 5 years. He tests a sample of 40 bulbs and finds the average time to be 4.8 years.

- (a) State null and alternate hypothesis
- (b) At a 2% significance level, is there enough evidence to support the idea that the warranty should be revised?

$$\text{Ans} \quad \mu = 5 \quad \sigma = 0.50 \quad n = 40 \quad \bar{x} = 4.8$$

- a) Null Hypothesis  $H_0: \mu = 5$   
 Alternate Hypothesis  $H_1: \mu < 5$  {1 tail test}

### 5) Decision Boundary



### c) Z-test

$$Z_d = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{4.8 - 5}{0.50 / \sqrt{40}}$$

$$= -2.53$$

Area under curve with Z score  $-2.53 = 0.0570$ .

$$P\text{-Value} = 0.0570 \quad \alpha = 0.02$$

Compare P-Value with  $\alpha$

$$0.0570 < 0.02 \Rightarrow \text{False}$$

We accept the Null Hypothesis

,

We Fail to Reject the Null Hypothesis.

## Student t distribution

In Z stats when we perform any analysis using Z-score  
we require  $\sigma$  (population standard deviation)  $\rightarrow$  is already known

How do we perform any analysis when we don't know  
the population standard deviation?



Student's t distribution

t stats

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

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

$s$  = Sample standard deviation

Z table

t table  $\Rightarrow$  t test

Degree of freedom

$$dof = n - 1 = 3 - 1 = 2$$

3 people



T-stats  $\div$  J test  $\rightarrow$  One Sample t-test.

- ① In the population the average IQ is 100. A team of researchers 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? CI = 95%  $\alpha = 0.05$

Ans)  $\mu = 100$   $n = 30$   $\bar{x} = 140$   $s = 20$  CI = 95%  $\alpha = 0.05$

① Null Hypothesis  $H_0 \div \mu = 100$

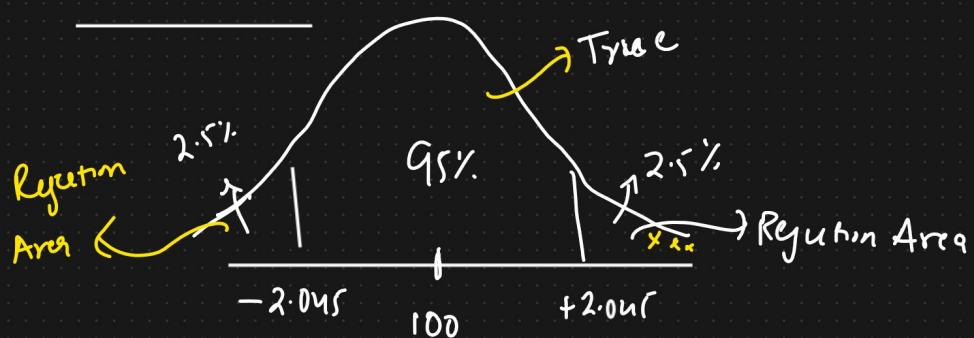
Alternate hypothesis  $H_1 \div \mu \neq 100$  {2 Tail Test}

②  $\alpha = 0.05$

③ Degree of freedom

$$df = n - 1 = 30 - 1 = 29.$$

④ Decision Rule



if t test is less than  $-2.045$  or greater than  $2.045$ , reject the null hypothesis

## ⑤ Calculate Test statistics

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{140 - 100}{20/\sqrt{30}} = \frac{40}{3.65}$$

$$t = 10.96$$

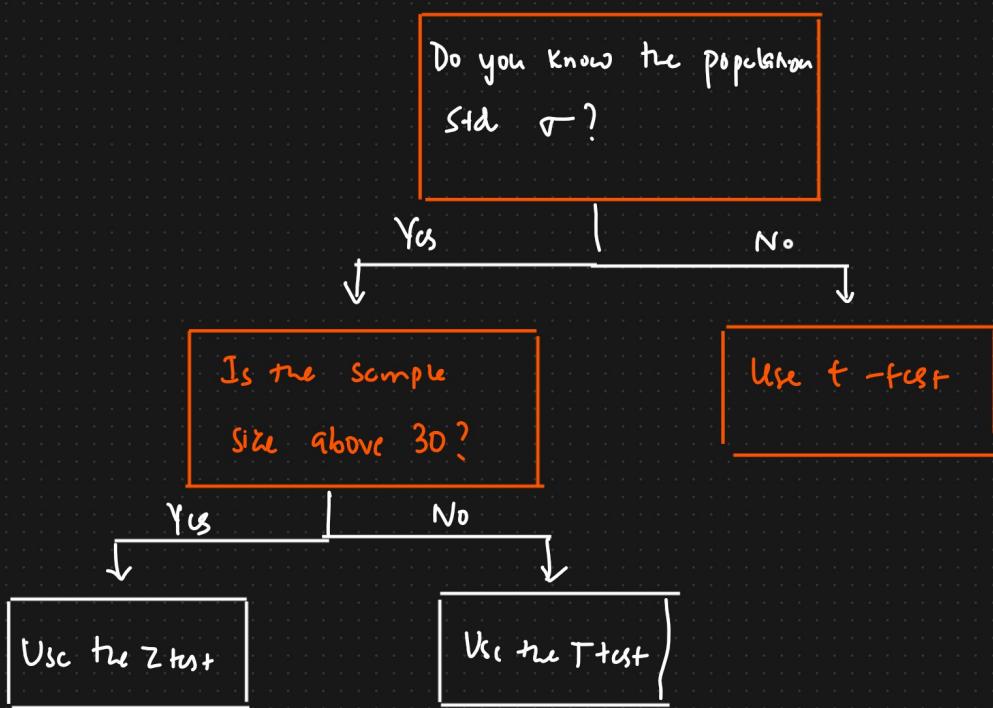
Since

$t = 10.96 > 2.045$  {Reject the Null Hypothesis}.

Conclusion : Medication used has affected the Intelligence

Medication has increased the Intelligence

## When To Use T-test Vs Z-test



## Type 1 and Type 2 Errors

Reality : Null Hypothesis is True or Null Hypothesis is False

Decision : Null Hypothesis is True or Null Hypothesis is False

Outcome 1 : We reject the Null Hypothesis when in reality  
it is false → Good

Outcome 2 : We reject the Null Hypothesis when in reality  
it is True → Type 1 Error

Outcome 3 : We retain the Null Hypothesis, when in reality  
it is False → Type 2 Error

Outcome 4 : We retain the Null Hypothesis when in  
reality it is True → Good

# Bayes Statistics (Bayes Theorem)

Bayesian statistics is an approach to data analysis and parameter estimation based on Bayes' theorem.

## Bayes' Theorem

Probability  $\begin{cases} \rightarrow \text{Independent Events} \\ \rightarrow \text{Dependent Events} \end{cases}$

### ① Independent Events

Eg: Rolling a dice

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

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

Tossing a coin

$$Pr(H) = 0.5 \quad Pr(T) = 0.5$$

### ② Dependent Event

$$\begin{array}{c} \text{Red} \rightarrow Pr(R) = \frac{2}{5} \xrightarrow{\text{Yellow}} Pr(Y) = \frac{3}{4} \\ \boxed{000} \\ 00 \end{array}$$

$$Pr(R \text{ and } Y) = Pr(R) * \boxed{Pr(Y|R)}$$

$$= \frac{2}{5} * \frac{3}{4} = \frac{6}{20} //$$

$$Pr(A \text{ and } B) = Pr(B \text{ and } A)$$

$$Pr(A) * Pr(B|A) = Pr(B) * Pr(A|B)$$

$$\boxed{Pr(B/A) = \frac{Pr(B) * Pr(A|B)}{Pr(A)}} \Rightarrow \text{Bayes' theorem}$$



$$P_{\delta}(A|B) = \frac{P_{\gamma}(A) * P_{\gamma}(B|A)}{P_{\gamma}(B)}$$

$A, B$  = events

$P_{\gamma}(A|B)$  = Probability of  $A$  given  $B$  is true

$P_{\delta}(B|A)$  = " " " $B$ " " $A$  is true

$P_{\gamma}(A), P_{\gamma}(B)$  = Independent probabilities of  $A$  and  $B$

<u>DATASET</u>		$\uparrow$ Independent	$\uparrow$ O/p / dependent
Size of Movie	No. of Rooms	location	Price
$x_1$	$x_2$	$x_3$	$y$

$$P_{\gamma}(y|x_1, x_2, x_3) = \frac{P_{\gamma}(y) * P_{\gamma}(x_1, x_2, x_3|y)}{P_{\gamma}(x_1, x_2, x_3)}$$

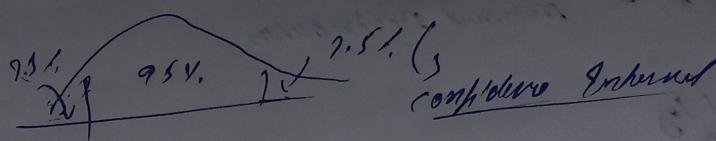


Bayes' Theorem

Section 23: Informed statistics

23. Confidence interval and margin of error.

$$c.t = 95\%$$



Painst Estimato

$$\boxed{\pi} \rightarrow \boxed{u}$$

Confidence Interval

Painst estimate  $\pm$  Margin of Error

$$Z \text{ test} \Rightarrow \pi \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

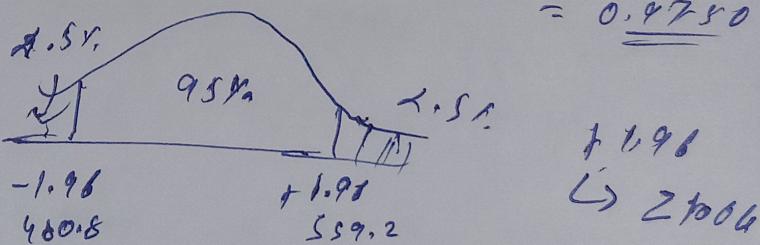
Ex: On the verbal section of CAT exam, the standard deviation is known to be 100. A sample of 20 test takers has a mean of 320. Construct 95% C.T. about the mean

$$\sigma = 0.08$$

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

$$= 0.025$$

$$= 0.9750$$



$$\text{Lower C.T.} = 320 - (1.96) \times \frac{100}{\sqrt{20}} = 480.8$$

$$\text{Higher C.T.} = 320 + 1.96 \times \frac{100}{\sqrt{20}} = 559.2$$

Ans: From 95% confidence about the mean CAT score is between 480.8 and 559.2

## CHI-SQUARE TEST

The Chi-square test for goodness of fit test claims about population proportions.

It is a non parametric test that is performed on categorical [ordinal and nominal] data.

There is a population of male who like different colors like

	<u>Theory</u>	<u>Sample</u>	④ Goodness of fit test
Yellow Bike	1/3	22	
Red Bike	1/3	17	
Orange Bike	1/3	39	
	1/3	78	↳ Observed
Theory Categorical distribution			categorical distribution

## Goodness of fit test

In a science class of 75 students, 11 are left handed, Does this class fit the theory that 12% of people are left handed.

1)

	<u>O</u>	<u>E</u>	$\frac{12}{100} \times 75 = 9$
left handed	11	9	
right handed	64	66	

## CHI-SQUARE for Goodness of fit

Chi-square for Goodness of fit ( $\chi^2$  for goodness of fit)

In 2010 census of the city - the weight of the individuals in a small city were found to be the following

$R_{50\text{kg}}$	$50-75$	$> 75$
20%	30%	50%

In 2020, weight of  $n=500$  individuals were sampled.

Below are the results.

$\leq 50$	$50-75$	$> 75$
140	160	200

Using  $\alpha = 0.05$ , would you conclude the population differences of weight has changes in the last 10 years?

A.v.  
=

Expected =

$\leq 50$	$50-75$	$> 75$
$0.2 \times 500$ $= 100$	$0.3 \times 500$ $= 150$	$0.5 \times 500$ $= 250$

- ② Null Hypothesis:  $H_0$ : The data meets the expectation.  
 Alternative Hypo:  $H_1$ : The data does not meet the expectation.

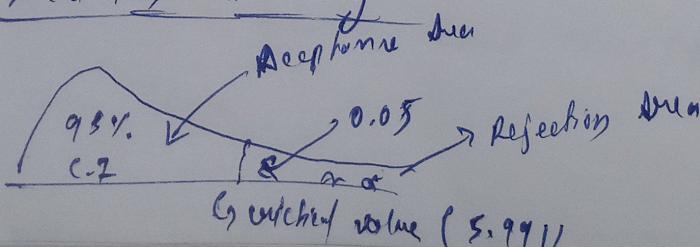
③  $\alpha = 0.05$ ,  $P.T = 95\%$  (confidence Interval)

- ④ Degrees of freedom

$$df = R - L = 3 - 1 = 2$$

$$c.v = df \times \alpha$$

- ⑤ Decision Boundary



If  $\chi^2$  is greater than 5.99, Reject No.  
else  
we will fail to reject the null hypothesis.

③ Calculate the various test statistics

$$\chi^2 = \sum \frac{(O - E)^2}{E} > \text{expected}$$

2020  
 $n=500$   
observed

<50	50-75	>75
140	160	200

Expected

<50	50-75	>75
0.2 \times 500	0.3 \times 500	0.5 \times 500
=100	=150	=250

$$= \frac{(140 - 100)^2}{100} + \frac{(160 - 150)^2}{150} + \frac{(200 - 250)^2}{250} = 28.8$$

$$= \frac{1600}{100} + \frac{100}{150} + \frac{2500}{250} = 10 + 0.66 + 10 = 20.66$$

$$\chi^2 = 20.66$$

If  $\chi^2$

$20.66 > 5.99$ , Reject No.

Answer

The weight of 2020 population are different than those expected in 2010 paper.

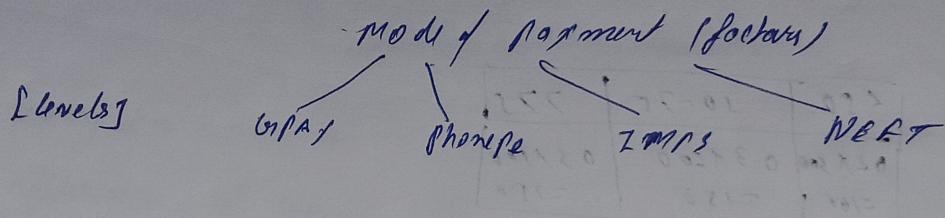
## Analyses of variance (ANOVA)

Def'n: ANOVA is a statistical method used to compare the means of 2 or more groups.

### ANOVA

- ① factors (variables)
- ② Levels,

e.g.: medicine (factor)  
 [Dosage]      5mg    10mg    15mg → levels



## Analyses of variance (ANOVA)

### Assumptions in ANOVA

- ① Normality of sampling distribution of mean  
 The distribution of sample mean is normally distributed.
- ② Absence of outliers  
 Outlying score need to be removed from the data set.
- ③ Homogeneity of variance

Population variance in different level of each independent variable are equal.

$$\sigma^2_1 = \sigma^2_2 = \sigma^2_3$$

① samples are independent and random.

## # Analysis of Variance (ANOVA)

### Types of ANOVA (3 types)

① One way ANOVA: One factor with at least 2 levels, these levels are independent.

E.g Doctor wants to test a new medication to decrease headache. They split the participation in 3 conditions (10mg, 20mg, 30mg). Doctor ask the participants to note the headache P1-P10.

Medication  $\rightarrow$  factor  
10 mg      20 mg      30 mg

② Repeated measure ANOVA:- One factor with at least 2 levels, levels are dependent.

Running  $\rightarrow$  factor  
levels  $\rightarrow$  Dog 1      Dog 2      Dog 3  
                  8                5                7  
                  7                4                9

③ Factorial ANOVA! Two or more factors (each of which with at least) 2 levels, levels can be independent and dependent

Gender  
↓  
factors

Male
Female

factor  
Running  $\rightarrow$  factor  
Dog 1 (level)      Dog 2      Dog 3  
8                5                9  
9                4                3  
2                6

## Analysis of Variance (ANOVA)

Hypothesis testing in ANOVA (partitioning of variance in ANOVA)

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

Alternative Hypothesis  $H_1$ : At least one of the sample means is not equal.  
 $\underline{(\mu_1 \neq \mu_2 \neq \mu_3)}$

### Test statistics

$$F = \frac{\text{variance b/w samples}}{\text{variance within samples}}$$

$$\mu_1 = \bar{x}_1 = \bar{x}_2 = \bar{x}_3$$

		variance b/w samples			$H_1: \text{At least one sample mean is not equal.}$
		$\bar{x}_1$	$\bar{x}_2$	$\bar{x}_3$	
variance within samples	1	1	6	5	
	2	7	6		
	3	3	3		
	4	2	2		
	5	1	4		
	6			$\bar{x}_3 = 4$	
		$\bar{x}_1 = 3$	$T_2 = 19/4$		

## One way ANOVA

One factor with at least 2 levels, levels are independent.

- ① Doctor wants to test a new medication which reduces headache. They splits the participant into 3 conditions [15mg, 30mg, 45mg]. Doctor ask the doctor ask the 10 people to headache b/w 1-10.

Are there any diff. b/w the 3 conditions  
 $\alpha = 0.052$

15 mg	30 mg	45 mg
9	7	4
8	6	3
7	6	2
8	7	3
8	8	5
9	7	3
8	6	2

① Define Null and Alternative Hypothesis:

$$H_0: \mu_{15} = \mu_{30} = \mu_{45}$$

$H_1$ : not all  $\mu$  are equal.

② Significance  $\alpha = 0.05$

③ Calculate Degree of freedom

$$\begin{matrix} N=21 & q=3 & n=7 \\ \text{sample} & \text{categories} & \end{matrix}$$

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

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

$$df_{\text{total}} = N-1 = 20$$

df<sub>1</sub>, df<sub>2</sub>

(2, 18)

↓

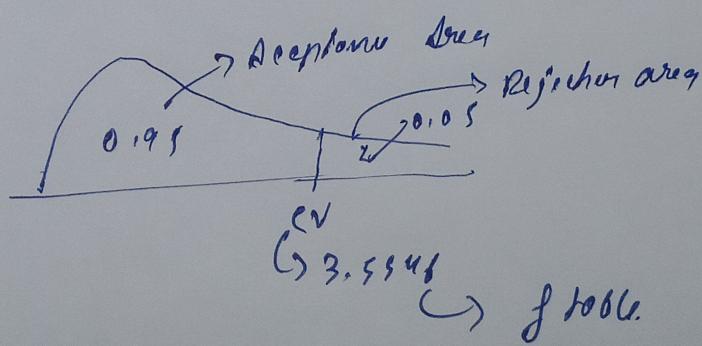
F table

$$\alpha = 0.05$$

↓

critical value

④ Decision Boundary



Decision Rule

If F is greater than 3.5546, reject the Null hypothesis

⑤ Calculate F test statistic

sum of squares (SS)

MS (Mean square)

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

	SS	df	MS	F
Between	98.67	2	49.33	
Within	10.29	18	0.57	
Total	108.96	20		

$$① SS_{\text{between}} = \frac{\sum (x_i - \bar{x})^2}{n} = \frac{7^2}{7} = 7$$

$$SS_{\text{within}} = 418.17 + 18.18 + 9.18 = 57$$

$$SS_{\text{total}} = 7 + 8 + 8 + 2 + 18 + 7 + 6 = 42$$

$$SS_{\text{between}} = 413.72 + 37.413.72 = 21$$

$$= \frac{57^2 + 47^2 + 21^2}{7} = \frac{57^2 + 47^2 + 21^2}{21}$$

$$= \underline{\underline{198.67}}$$

$$\textcircled{Q} \quad S_{\text{within}}^2 = \bar{y}^2 - \frac{\sum (\bar{y}_i)^2}{n}$$

$$\bar{y}^2 = 9^2 + 1^2 + 2^2 + 8^2 + 8^2 + \dots$$

$$= 853$$

$$= 853 - \frac{[52^2 + 47^2 + 21^2]}{7}$$

$$= 10.29$$

$$f_{\text{test}} = \frac{M_s \text{ Between}}{M_s \text{ Within}}$$

$$F = \frac{\text{variance b/w samples}}{\text{variance b/w within samples}}$$

$$F' = \frac{49.39}{0.54} = 86.58$$

If F is greater than 3.5548, Reject H<sub>0</sub>.

86.58 > 3.5548      Reject H<sub>0</sub>.