

We can not believe in it.

0.05

Table

(9-9)

)

Central limit Theorem for (Population) proportion %

Distribution of proportions follows approximately normal for a large sample distribution with  $P$  (The population proportion), and standard deviation  $\sqrt{\frac{P(1-P)}{n}}$

Hypothesis Testing for Population proportion %:

→ Z-Test for proportion according to Central limit Theorem of proportions the Sampling distribution of proportion  $\bar{P}$  for large sample follows an approximate normal distribution

$$Z = \frac{\bar{P} - P}{\sqrt{\frac{P(1-P)}{n}}}$$

To calculate standard deviation  $\sqrt{P(1-P)}$

We need knowledge of 'p', however we can use the value of  $\hat{p}$  estimated from large samples.

One of the thumb rules used is.

→ That the value of  $n \times \hat{p} \times (1-\hat{p}) \geq 10$ , to use

$$z = \frac{(\hat{p} - p)}{\sqrt{\frac{p(1-p)}{n}}}$$

- ① According to a study exactly 62% gift card purchased from e-commerce were never used. the Manager at e-commerce company wanted to test whether this claim is true. collected data 250 gift card purchases and found the 22 gift card's for not used till it's expiry date.

$$n = 250$$

$$H_0: p \leq 0.12$$

$$H_1: p > 0.12$$

$$p = \frac{12}{100}$$

$$p = 0.12$$

$$n \times \hat{p} \times (1-\hat{p}) \geq 10$$

$$250 \times 0.088 \times (1 - 0.088) \geq 10$$

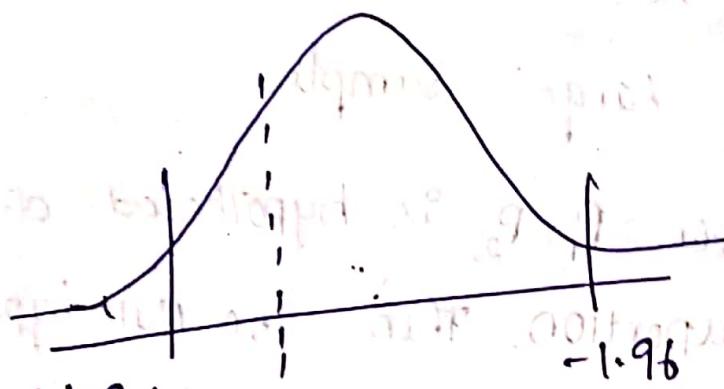
$$= 22 \times 0.912 = 20.064$$

$$z = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} = \frac{0.088 - 0.12}{\sqrt{\frac{0.12(1-0.12)}{250}}} = \frac{0.88 + 0.02}{\sqrt{0.88 + 0.02}}$$

$$= \frac{0.88 - 0.12}{0.02055} = \frac{0.76}{0.02055}$$

$$= \frac{-0.032}{0.02055} = -1.557$$

It is two tailed Test  $\alpha = 0.05$   $Z = 1.96$



$H_0$  is accepted. can not be rejected.

confidence Interval  $\bar{x} \pm z \sigma / \sqrt{n}$

$$\bar{p} \pm 1.96 \sqrt{\frac{p(1-p)}{n}}$$

$$0.088 \pm 1.96 \left( \sqrt{\frac{0.12(1-0.12)}{250}} \right)$$

$$[0.0528, 0.1231]$$

Unbiased procedure not responsive problem with (1).

Impact of significant effect size on the test statistic, non detection of significant difference as long as a large effect size is present in the population.

18/09/2018  
 Hypothesis Test for difference Population proportion, TWO Sample Z-Test of Proportion  
 → When the proportions estimated from Large sample then sampling distribution of proportions follow a normal distribution according to Central Limit Theorem.

Let  $\bar{P}_1$  &  $\bar{P}_2$  be estimated values of proportion from Large sample:

The Difference  $P_1 - P_2$  is hypothesized difference in population proportion. Then the Null hypothesis is  $P_1 = P_2$  (i.e.  $H_0: P_1 = P_2$ ) then the Test

Statistic is

$$Z = \frac{(\bar{P}_1 - \bar{P}_2) - 0}{\sqrt{\bar{P}(1-\bar{P}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

Where  $\bar{P} = \frac{n_1 \bar{P}_1 + n_2 \bar{P}_2}{n_1 + n_2}$  is the pooled estimate for proportion.

- ① The Marketing Manager for a company believes that non-affluent customers are sensitive to discount compare to affluent customers to validate this hypothesis discount coupons are send to non-affluent and affluent customers and the data is provided below.

Group	Sample Size	No of customers using coupons	Estimated proportion
Non-Affluent	500	145	0.29
Affluent	300	42	0.14

use an appropriate hypothesis test to whether  
 @ there is any difference in proportion of  
 customer who use discount coupons at  
 $\alpha = 0.05$ .

ns: The Null + Alternative hypothesis

$$H_0: P_1 = P_2$$

$$H_1: P_1 \neq P_2$$

$$\bar{P}_1 = 0.29, \quad \bar{P}_2 = 0.14.$$

$$\text{pooled proportion } \bar{P} = \frac{n_1 \bar{P}_1 + n_2 \bar{P}_2}{n_1 + n_2}$$

$$= \frac{500 \times 0.29 + 300 \times 0.14}{500 + 300}$$

$$= \frac{145 + 42}{800} = \frac{187}{800}$$

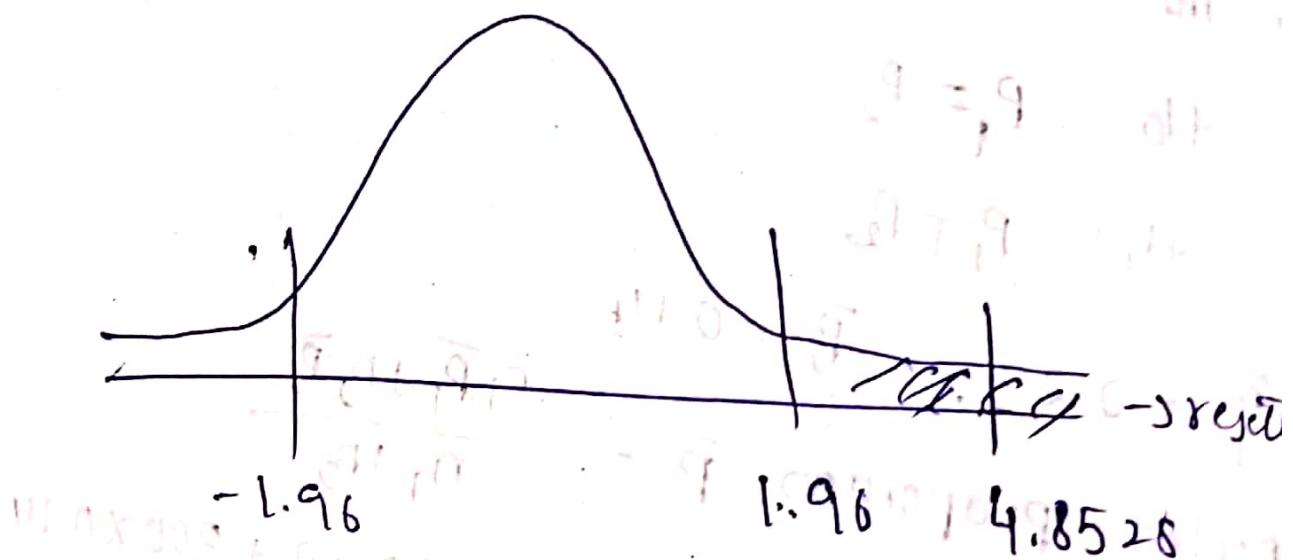
$$= 0.2338.$$

$$Z = \frac{(\bar{P}_1 - \bar{P}_2) - 0}{\sqrt{\bar{P}(1-\bar{P})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

$$= \frac{0.29 - 0.14}{0.2338\sqrt{1 - 0.233}}$$

$$\begin{aligned}
 &= \frac{0.29 - 0.14}{0.2338(1 - 0.2338)\left(\frac{1}{500} + \frac{1}{300}\right)} \\
 &= \frac{0.5}{0.179(0.002 + 0.0033)} \\
 &= \frac{0.15}{0.030091} = 4.8528
 \end{aligned}$$

at 0.05 two tail z critical value 1.96  
 so we can reject Null hypothesis



23/09/2019

F-Distribution = (Fisher's distribution)

⇒ Ratio of two  $\chi^2$  distributions is follow

F-Distribution

⇒ F-Distribution is a ratio of two  $\chi^2$

distributions

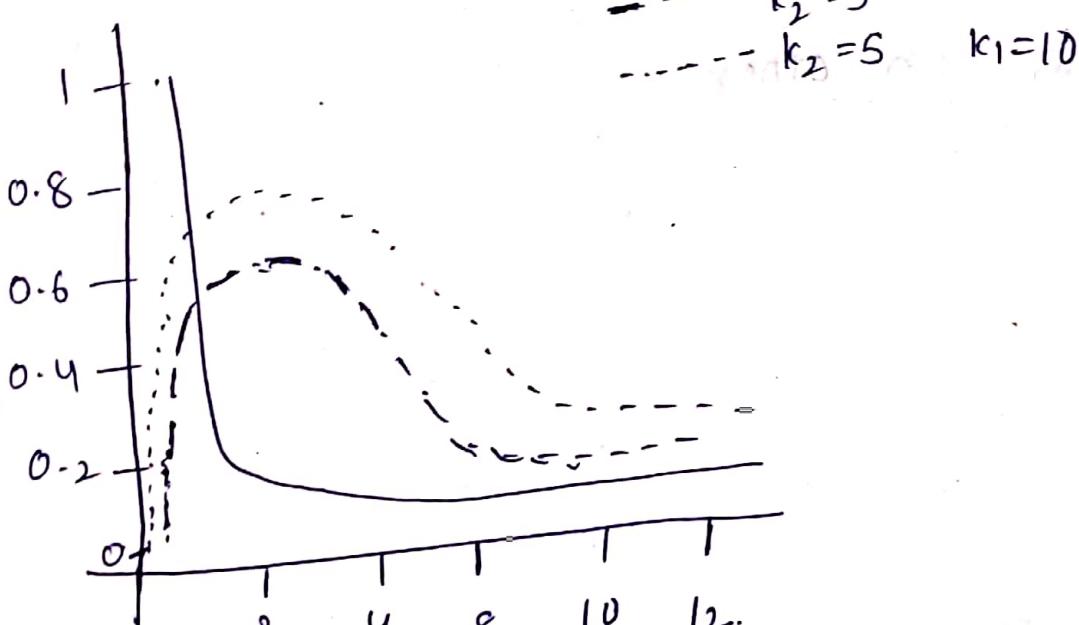
Let  $Q_1$  and  $Q_2$  be two independent  $\chi^2$

distributions with  $k_1$  &  $k_2$  degree of freedom

respectively. Then random variable 'F' defined

as  $F = \frac{Q_1/k_1}{Q_2/k_2}$  is an 'F' Distribution.

$$f(x) = \frac{T\left(\frac{k_1+k_2}{2}\right)\left(\frac{k_1}{k_2}\right)^{k_1/2} x^{\frac{k_1}{2}-1}}{T\left(\frac{k_1}{2}\right)+\left(\frac{k_2}{2}\right)} \left(1+\frac{k_1-x}{k_2}\right)^{\frac{k_1+k_2}{2}}$$



properties 'F' distribution

1) Mean of 'F' distribution is  $\frac{k_2}{(k_2-2)}$   $k_2 > 2$

2) Standard deviation of 'F' distribution

is 
$$\sqrt{\frac{2k_2^2(k_1+k_2-2)}{k_1(k_2-2)^2(k_2-4)}}$$
 for  $k_2 > 4$

$k_1, k_2$  are degree of freedom

3) 'F' distribution is non-symmetrical and the shape of the distribution depends on values  $k_1$  and  $k_2$

4) 'F' distribution is used in Analysis of variance to Test the mean values of multiple groups

\* with Two - Tail F Test we just want to If the variance are not equal to each other

① Conduct a Two-tail 'F' Test and the following samples.

sample 1

$$\text{Variance} = 109.63$$

$$n_1 = 41$$

sample 2

$$\text{Variance} = 65.99$$

$$n_2 = 21$$

$$H_0: \sigma_1^2 = \sigma_2^2$$

$$H_1: \sigma_1^2 \neq \sigma_2^2$$

Calculate 'F' statistic

$$\text{which } F = \frac{\sigma_1^2}{\sigma_2^2} = \frac{\sum (x - \bar{x}_1)^2 / (n-1)}{\sum (x - \bar{x}_2)^2 / (n-1)}$$

It will follow 'F' distribution.

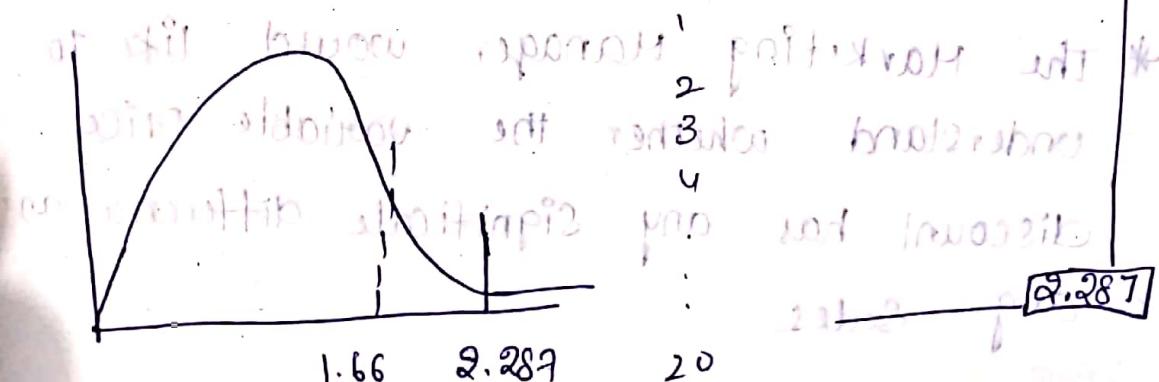
$$F = \frac{109.63}{65.99} = 1.66.$$

Assume  $\alpha = 0.05$  and degree freedom  $df_1 = 40$

$df_2 = 20$

critical value of  $F$  stat  $df_1 = 1.27300$

$df_2 = 2.287$



We can not reject null hypothesis since 'F' statistic value less than critical value so we can not reject null hypothesis.

Since two variance not significant difference

24/09/2019

## Analyses of variance (ANOVA)

→ In many situations we need to conduct a hypothesis to compare mean value of samples created by using a factor in two groups (sample).

Ex: A marketer like to understand the impact of '3' different discount values such (0%, 10%, 20%) on the average sales. When we have to compare the impact of a factor on mean of more than two groups. Since the hypothesis test such two sample T-test not a ideal approach since they can result in incorrect Type I and Type II error. We can use the analyses of variance to understand the difference in population mean among more than two populations.

\* The Marketing Manager would like to understand whether the variable price discount has any significant difference on avg sales.

\* In ANOVA our objective is to verify whether variance due to treatment is different from variance due to randomness.

One way Analysis of Variance:  
→ One way ANOVA appropriate under following condition  
\* We would like to study Impact of single Treatment (factor) at different Level on continuous Responses variable (outcome).

For Example The variable Price discount is the treatment factor and 0%, 10% and 20% price discounts are different levels (3 levels) different levels of discounts are likely to have a varying impact on sales of the product where Sales is outcome variable.

\* In each group the population response variable follows a normal distribution and the sample chosen using random sampling.

\* Population variance for different groups are assumed to be same. i.e. variability in the response variable values with in different groups is same.

## Setting up Analysis of Variance:

\* Assume that we would like to study the impact of a factor (discount) with  $k$  levels on ~~continuous~~ continuous variable (such as Sales quantity). Then the null and alternate hypothesis for one way ANOVA

$$H_0 : \mu_0 = \mu_1 = \mu_2 = \dots = \mu_k$$

$$H_a : \text{Not all } \mu_i \text{ are equal}$$

- ① Using the following data perform a one way Analysis of variance using  $\alpha = 0.05$

Group 1      Group 2      Group 3

51	23	56
45	43	76
33	23	74
45	43	87
67	45	56

$$\text{Group 1} = 51 + 45 + 33 + 45 + 67$$

$$\text{Mean} = \frac{51 + 45 + 33 + 45 + 67}{5} = 48.2$$

$$\text{Group 2} = \frac{23 + 43 + 23 + 43 + 45}{5} = 35.4$$

$$\text{Group 3} = \frac{56 + 76 + 74 + 87 + 56}{5} = 69.8$$

GROUP <sub>1</sub>		VARIANCE	DEVIATION	SQUARE DEVIATION
51	48.2		51 - 48.2 = 2.8	(2.8) <sup>2</sup> = 7.84
45	48.2		-3.2	10.24
33	48.2		-15.2	231.04
45	48.2		-3.2	10.24
67	48.2		18.2	354.44

GROUP <sub>2</sub>		DEVIATION	SQUARE.
23	35.4	-12.4	153.76
43	35.4	7.6	57.76
23	35.4	-12.4	153.76
43	35.4	7.6	57.76
45	35.4	9.6	92.16

GROUP <sub>3</sub>		DEVIATION	SQUARE.
56	69.8	-13.8	190.44
76	69.8	6.2	38.44
74	69.8	4.2	17.64
81	69.8	17.2	295.84
56	69.8	-13.8	190.44

$$\text{Sum of square deviation} = \frac{-7.84 + 10.24 + 231.04 + 10.24 + 354.44}{5}$$

$$(G_1) = 612.84$$

$$= \frac{153.76 + 57.76 + 153.76 + 57.76 + 92.16}{5}$$

$$= 515.211$$

$$= \frac{190.44 + 38.44 + 17.64 + 295.84 + 190.44}{5}$$

$$(G_3) = 732.811$$

$$\text{var}_1 = 612.8 / 5-1 = \frac{612.8}{4} = 153.2$$

$$\text{var}_2 = \frac{515.2}{5-1} = \frac{515.2}{4} = 128.8$$

$$\text{var}_3 = \frac{732.8}{4} = 183.2$$

$$\text{Mean Square Error} (\text{MS}_{\text{Error}}) = \frac{153.2 + 128.8 + 183.2}{3}$$

$$= 155.07.$$

25/09/2019

$$(\text{Grand } \bar{x}) (\bar{x}_{\text{Grand}}) = \frac{48.2 + 35.4 + 69.8}{3} = 51.13$$

Again find variances for group means

Group Mean Grand Mean

$$48.2 - 51.13$$

$$35.4 - 51.13$$

$$69.8 - 51.13$$

$$48.2 - 51.13$$

$$\frac{\sum (x - \bar{x})^2}{n-1}$$

$$= \frac{(48.2 - 51.13)^2 + (35.4 - 51.13)^2 + (69.8 - 51.13)^2}{3-1}$$

$$= \frac{604.58}{2} = 302.29$$

~~MSB~~ (Mean Square Variation with in Group)  $\boxed{\text{MSW}}$

MSB (Mean Square Variation between groups)

$$= 302.29 \times 5$$

$$= 1511.45$$

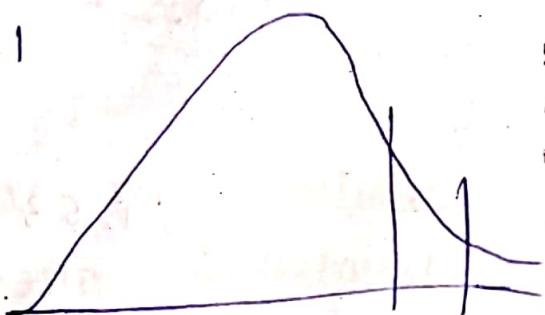
$$\begin{aligned} F\text{-Statistic} &= \frac{\text{MSB}}{\text{MSW}} \\ &= \frac{1511.45}{155.07} \\ &= \underline{\underline{9.75}} \end{aligned}$$

$$\begin{aligned} \text{Degree of freedom} &= 15 - 3 \\ \text{within group} &= 12 \end{aligned}$$

df	1	2
1		
2		
3		
:		
12		

$$\begin{aligned} \text{Degree of freedom} &= 3 - 1 \\ \text{across group} &= 2 \end{aligned}$$

$$df(2, 12) = 3.89$$



3.89 9.75

reject null hypothesis

With in the following summary data perform  
a one way analysis of variance using

n	mean	std.dev
30 G <sub>1</sub>	50.26	10.45
30 G <sub>2</sub>	45.32	12.76
30 G <sub>3</sub>	53.67	11.47

$$\text{variance}_1 = 10.45^2 = 109.2$$

$$\text{var}_2 = 12.76^2 = 162.82$$

$$\text{var}_3 = 11.47^2 = 131.56$$

$$F = \frac{MSB}{MSW}$$

$$MSW = \frac{109.2 + 162.82 + 131.56}{3}$$

$$= 134.53.$$

~~MSB~~

$$\text{Grand mean} = \frac{50.26 + 45.32 + 53.67}{3}$$

$$= 49.75$$

Group mean

$$= \frac{(50.26 - 49.75)^2 + (45.32 - 49.75)^2 + (53.67 - 49.75)^2}{3-1}$$

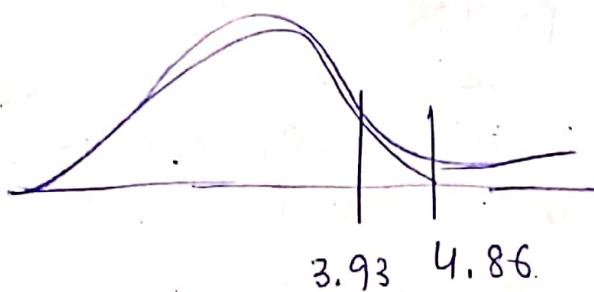
$$MSB = \frac{35.25}{2} = 17.62$$

$$\Rightarrow MSB = 17.62 \times 30$$
$$= 528.75$$

$$F_{\text{statistic}} = \frac{MSW}{MSR}$$

$$= \frac{528.75}{134.53} = 3.93$$

$$F_{\text{critical}}(2, 87) = 4.86.$$



All Group means are equal. So we can not reject null hypothesis

A clinical Psychologist has run a between subject experiment comparing two treatments for depression (formative behavioral therapy and client center therapy) against a control condition. Subjects for

Randomly assign the experimental conditions after 12 weeks. The subjects differ in

Score Measure using CESD scale differentiation scale that data are summarized as.

	n	mean	sd.dev
Control	40	21.4	4.5
CBT	40	16.9	5.5
CCT	40	19.1	5.8

Using a one way ANOVA with  $\alpha = 0.05$

$$M_{GW} = \frac{20.25 + 30.25 + 33.64}{3} = 28.05$$

Grand mean = 19.13.

$$\text{Group Mean} = \frac{(21.4 - 19.13)^2 + (16.9 - 19.13)^2 + (19.1 - 19.13)^2}{3}$$

$$= \frac{5.15 + 4.97 + 0.00}{3}$$

$$MSB = \frac{10.12}{2}$$

$$MSB = \frac{10.12 \times 40}{2} = 202.4$$

$$F_1 = \frac{MSB}{MSW} \quad F_{(1)} = \frac{MSB}{MSW}$$

$$F_{(1)} = \frac{28.05}{202.4} \quad F_{(1)} = \frac{202.4}{28.05}$$

$$F_{\text{critical}} (2, 11) = 4.79 \quad F_{(1)} = 7.22$$

We can reject null hypothesis.

# Covariance analysis

26/09/19

Two variance indicate the direction of linear relationship between variables

## Correlation

Correlation measure both strength and direction of the linear relationship between two variables.

Covariance of two variables  $x, y$ ,  $\text{cov}(x, y)$  can be represented as  $\text{cov}(x, y)$

If  $E(x)$  is expected value or mean of a sample  $x$ . then  $\bar{x}$

$$\text{cov}(x, y) = E[(x - \bar{x})(y - \bar{y})]$$

Correlation Coefficient (Pearson correlation coefficient)  
⇒ Pearson product moment

is used for measuring the strength and direction of linear relationship between two continuous random variables  $x, y$ .

Let  $x_i$  be different value of variable  $x$  and  $y_i$  be different value of variable  $y$  then Pearson correlation coefficient

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{n \sigma_x^2 \sigma_y^2}}$$

When we try to measure how change in variable [Y] is related to changing another variable (X) one of the issues that we need consider is the measurement scale and unit of measurement of two variables.

For example age is measured in years and call duration is measured in seconds.

The range of two variables can be different, thus we need standardised variables which can be used for measuring correlation between two variables.

$$z_x = \left( \frac{x_i - \bar{x}}{\sigma_x} \right) \quad z_y = \left( \frac{y_i - \bar{y}}{\sigma_y} \right)$$

$$r = \sum_{i=1}^n \frac{z_x \cdot z_y}{n}$$

$$\gamma = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

$$\boxed{\gamma = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}}$$

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

### Properties:-

- ⇒ positive value indicate positive Correlation.
- ⇒ negative value indicate Negative Correlation.
- ⇒ Pearson correlation coefficient value may be even there is strong Non-relationship between  $x, y$ . Thus low correlation coefficient value can not be taken as evidence no relationship.

① The Avg share ratios of two companies over the first 3 months

$x$	$y$	$x_i - \bar{x}$	$y_i - \bar{y}$
274.58	219.50	274.58 - 284.286	-16.6
287.96	242.92	-9.706	3.68
290.35	245.90	3.68	6.82
<u><math>\bar{x} = 284.286</math></u>	<u><math>\bar{y} = 236.10</math></u>	<u>-6.064</u>	<u>9.82</u>

$$(x_i - \bar{x})^2 (y_i - \bar{y})^2 (x_i - \bar{x})(y_i - \bar{y})$$

$$94.08 \cdot 275.56 + 161.02$$

$$13.54 \cdot 46.51 + 25.09$$

$$36.84 \cdot 96.43 + 59.4$$

$$\underline{144.42} \quad \underline{418.1} \quad \underline{245.6}$$

$$r = \frac{245.6}{\sqrt{144.42} \times \sqrt{418.1}}$$

$$\frac{245.6}{12.01 \times 20.45}$$

$$r = \frac{245.6}{12.01 \times 20.45}$$

$$r = \frac{245.6}{245.60}$$

$$\therefore r = 1$$

Positively co-related and strongly

co-related

$\alpha/2$

### Simple Linear Regression:

Simple linear regression is a statistical for finding the existence of an association b/w the dependent variable and independent variable.

- Simple linear regression implies that there is only one independent variable in the order.

### Functional form of SLR:-

$$y = \beta_0 + \beta_1 x + \epsilon$$

\* For a dataset with  $n$

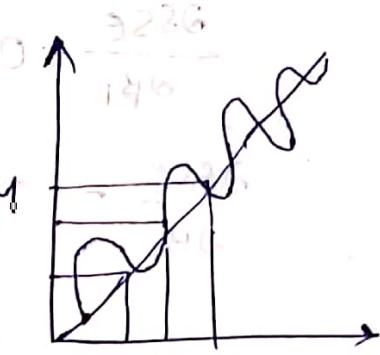
Observations  $(x_i, y_i)$  where

$$i = 1, 2, 3, \dots, n$$

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

where  $\epsilon$  - error

$$\epsilon = u_i - \beta_0 - \beta_1 x_i$$



\* Estimation of parameters using ordinary least squares:

$$\sum_{i=1}^n \varepsilon_i = \sum_{i=1}^n (y_i - B_0 - B_1 x_i)$$

- In ordinary least squares the objective is to find the optimal value of  $B_0$  and  $B_1$ , that will minimize the sum of squares error (SSE)

$$SSE = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - B_0 - B_1 x_i)^2$$

- To find the optimal values of  $B_0 + B_1$  that minimize SSE, we have to equate the partial derivatives of SSE with respect to  $B_0$  &  $B_1$  to 0.

$$\frac{\partial SSE}{\partial B_0} = 0$$

$$\text{and } \frac{\partial SSE}{\partial B_0} = - \sum_{i=1}^n 2(y_i - B_0 - B_1 x_i) = 0$$

$$B_0 = \bar{y} - B_1 \bar{x}$$

$$\frac{\partial SSE}{\partial B_1} = 0$$

$$\frac{\partial SSE}{\partial B_1} = - \sum_{i=1}^n 2x_i(y_i - B_0 - B_1 x_i) = 0$$

$$= -2 \sum_{i=1}^n (y_i x_i - B_0 x_i - B_1 x_i^2) = 0$$

Substitute  $B_0$  value in above eq

$$= -2 \sum_{i=1}^n (y_i x_i - \bar{y} \bar{x} + \beta_1 x_i \bar{x} - \beta_1 \bar{x} x_i)$$

$$= \sum_{i=1}^n (x_i y_i - \bar{y} \bar{x}) - \beta_1 \sum_{i=1}^n (x_i^2 - \bar{x} \bar{x}) = 0$$

$$\beta_1 = \frac{\sum_{i=1}^n (x_i y_i - \bar{y} \bar{x})}{\sum_{i=1}^n (x_i^2 - \bar{x} \bar{x})}$$

$$\beta_1 = \sum_{i=1}^n \frac{x_i (y_i - \bar{y})}{x_i (x_i - \bar{x})}$$

The above  $\beta_1$  can be simplified as

$$\text{since } \sum_{i=1}^n (\bar{x} \bar{y} - \bar{y} \bar{x}) = 0 \text{ & } \sum_{i=1}^n (\bar{x}^2 - \bar{x} \bar{x}) = 0$$

$$\beta_1 = \frac{\sum_{i=1}^n (x_i y_i - \bar{y} \bar{x}) + 0}{\sum_{i=1}^n (x_i^2 - \bar{x} \bar{x}) + 0}$$

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

$$\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

23/10/2019

Validation of simple Linear Regression:

1. Co-efficient as Determinant ( $R^2$ )
2. Hypothesis test for regression coefficient  $B_1$

3. ANOVA (MLR)

4. Residual Analysis

5. Outlier Analysis

Hypothesis testing for  $B_1$ :

Explained variance

$$R^2 = \frac{\text{Explained variance}}{\text{Total variation}}$$

$$y_i = \beta_0 + \beta_1 x + \epsilon_i$$

Variation in  $y_i$  = variation explained by model

SST = sum of square of total variation

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

SSR = SSR + SSE (Sum of squares of errors)

$$R^2 = \frac{\text{Explained variance}}{\text{Total variation}} = \frac{SSR}{SST} = \frac{SST - SSE}{SST}$$

$$R^2 = 1 - \frac{SSE}{SST}$$

\*  $R^2$  is the proportion of variation explained by the