

COLLECTION OF DATA

Before any statistical work can be done figure or data must be collected, it's very important for us to collect the right and correct data because any mistakes, error or bias in data collection will affect the result, decision taken and conclusion in research work.

DEMOGRAPHY.

Is a branch of the social sciences that deals with study of human population, their structure, change (through birth, death and migration) and their relationship with natural environment.

DEMOGRAPHIC INDICATORS.

These include population size, growth rate, birth rate, death rate, total fertility rate, life expectancy. Demographic changes affect all areas of human activities that is economic, social, cultural and political.

Definition of some demographic terminologies.

- (i) Crude birth rate: Is the annual number of live birth per 1000 people.
- (ii) Crude death rate: Is the annual number of death per 1000 people.
- (iii) Fertility rate: Is the annual number of live birth per 1000 women of child bearing age: 15 - 49 years.

- (5) Life expectancy: Is the number of years an individual at a given age is expected to live at present mortality level.
- (6) Total fertility rate: Is the number of live births per women completing her reproductive life if her child bearing at each age reflect current specific fertility rate.
- (7) The replacement level fertility: Is the average number of children a woman must have born to replace herself with daughter in the next generation.
- (8) The growth reproduction rate: Is the number of daughters who will be born to a woman completing her reproductive life at current age specific fertility rate.
- (9) The next reproduction rate: Is the expected number of daughters per new born prospective mother who may or may not survive to and through the age of bearing.
- (10) Stable population: Is a population that has a constant crude birth rate and death rate for a period of time and the percentage of people in every age limit remains constant.
- formulae:
- (i) Crude birth rate: $\frac{\text{Total number of birth}}{\text{Total number of population}} \times 1000$
 $(\text{per } 1000)$
- (ii) Crude death rate: $\frac{\text{Total number death}}{\text{Estimated population}} \times 1000$
 $(\text{per } 1000)$
- (iii) Age specific death rate (ASDR) = $\frac{\text{death rate}}{\text{Total Population}} \times 1000$
 $(\text{per } 1000)$

P.T.O

(i) Age specific birth rate: (ASBR)

$$= \frac{\text{Birth rate}}{\text{Total female female population}} \times 1000$$

Total female female population (per 1000).

(ii) Standard death rate: (SDR)

$$= \frac{\text{Total expected death}}{\text{Total standard population}} \times 1000$$

Total standard population (per 1000).

(iii) Expected death = Age specific death rate \times Standard population
(per 1000).

(iv) Standardized birth rate (SBR)

$$SBR = \frac{\text{Total expected birth}}{\text{Total standard female population}} \times 1000$$

Total standard female female

population (per 1000).

(v) Expected birth = Age specific birth rate \times Standard female female population
(per 1000).

Example: The data below shows the age group, estimated population, recorded death and standard population as recorded by a demographer.

Age group	Estimated population in million	Recorded death in thousands	standard population in millions
0 - 20	25	360	24
21 - 40	38	280	32
41 - 60	26	150	16
61 and above	8	650	8
Total	100	1446	80

Topic: final

- (a) (i) Crude death rate
- (ii) age specific death rate
- (c) (i) Standardized death rate

Solution:

(a) (i) Crude death rate (CDR) =

$$\frac{\text{Total number of deaths}}{\text{Estimated population}} \times 1000$$

$$= \frac{144,000}{100,000,000} \times 1000$$

$$= 144 \text{ per } 1000$$

(b) Age-specific death rates.

Age group	Estimated population	Record death	age specific death rate
0 - 20	28	360	$\frac{360}{28,000,000} = 12.857$
21 - 40	58	250	$\frac{250}{28,000,000} = 6.529$
41 - 60	26	180	$\frac{180}{28,000,000} = 6.923$
61 and above	85	650	$\frac{650}{28,000,000} = 23.214$
Total	100	1440	

Age group	Age-specific death rate	standardized	Expected deaths
0 - 20	12.857	24	308,568
21 - 40	6.529	32	210,528
41 - 60	6.923	16	110,768
61 and above	23.214	8	650,500
Sum		80	1,279,864

$$SADR = \frac{1,279,864}{80,000,000} \times 1000$$

$$30.8868$$

$$210,528$$

$$110,768$$

$$650,500$$

$$= 30.8868 \text{ per } 1000$$

$$= 16 \text{ per } 1000$$

Exer C5-C

- (1) The data below shows the age group, estimated population, estimated per 1000 female birth and standardized fertile female population as recorded by a statistician in a rural community in Nigeria.

Age group	Estimated Population (millions)	Estimated per 1000 female births	Estimated fertile female population	Standardized per 1000 female population
0 - 14	300	-	-	-
15 - 24	250	195	38	52
25 - 34	540	200	96	198
35 - 44	100	85	25	50
All infabiles	230	-	-	-
Total				

Required:

- (1) Calculate the crude birth rate
 (2) Compute the age specific and standardized birth.
- (3) The table below gives the population and the number of deaths in two towns A and B.

Age-group	Population (thousands)	Number of deaths (1000) A	Population (1000) B	Number of deaths (1000) B	Population (1000) B
0 - 2	2,300	150	3000	300	100
2 - 10	1,400	50	12,000	60	150
10 - 20	11,000	70	14,000	100	350
20 - 30	30,000	250	18,000	70	250
All infabiles	6,500	500	8,000	450	200



TOPIC Calculations

- (i) The crude death rate for the two towns.
- (ii) Its death rate for each group for each town.
- (iii) To standardize age group for each town.

REGRESSION: Is defined as the amount of change in one variable associated with a unit change in another variable. To be more specific for a variable y which depend on another variable x the independent variable and equation relating y on x is called a regression equation.

Regression equation is thus for a statistical tool or technique use to predict or explain the relationship between y and x . It's plain from the definition presented above the two types of variable are involved in regression analysis.

One variable changes in values and therefore induce a change in a unit value of another variables

i.e. a cause and effect relationship. x is called the independent variables explanatory variables, exogenous variables, control or regressor.

y is called the dependent variables, response, endogenous, uncontrol or the regress variable.

Regression can be simple or multiple linear or non-linear we shall discuss only the simple regression that is only two variables namely.

The Simple regression model

$$y_i = b_0 + b_1 x_i + \epsilon_i$$

Where y = dependent variable.

x_i = independent variables

b_0 = constant / intercept

b_1 = the regression coefficient
 ϵ_i = error terms.

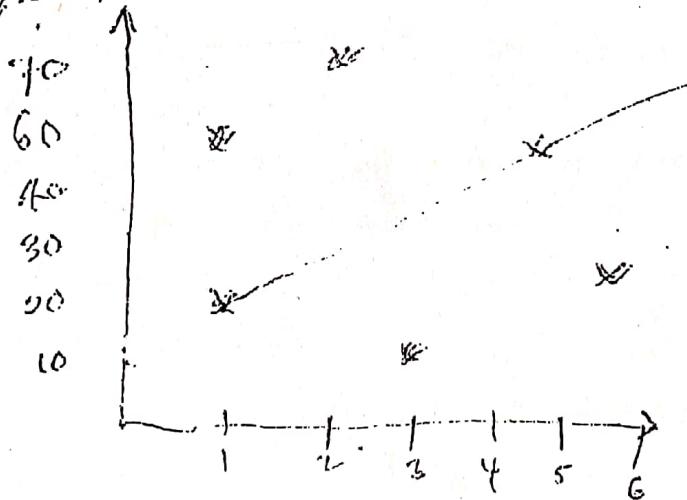
$$b_1 = \frac{n \sum xy - \bar{x} \bar{y}}{n \sum x^2 - (\bar{x})^2}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

where $\bar{y} = \frac{\sum y}{n}$ and $\bar{x} = \frac{\sum x}{n}$

BIVARIATE DATA: A bivariate data (x, y) is a set of values which appears in pairs it's a situation in which the values of y depend on the value of x .

SCATTERED Diagram: When the points (x, y) of a bivariate data are plotted in a rectangular coordinate system the resulting diagram is called a scatter diagram.



Simple CORRELATION.

Definition: It is a statistical techniques used to measure the degree, strength and direction of the relationship or association between two or more variables.
 The main purpose of Correlation is to find out how



Opic:

Strong or weak a relationship is or to find out IF a relationship is positive or negative, it's denoted by r_{xy} .

Types of Correlation.

- Pearson's product moment correlation co-efficient (PPM.C.C.)
- Spearman's Rank correlation coefficient.
- Kendall tau correlation coefficient.
- Partial and Multiple correlation formulae.

$$r_s = \frac{\sum xy}{\sqrt{(\sum x^2)(\sum y^2)}}$$

where $X = x - \bar{x}$ or $y = \frac{\sum xy}{\sum x^2}$

$$V_{xy} = \frac{N \sum xy - (\sum x)(\sum y)}{\sqrt{(N \sum x^2 - (\sum x)^2)(N \sum y^2 - (\sum y)^2)}}$$

If $r = 0.1$ — 0.49, a weak positive correlation.

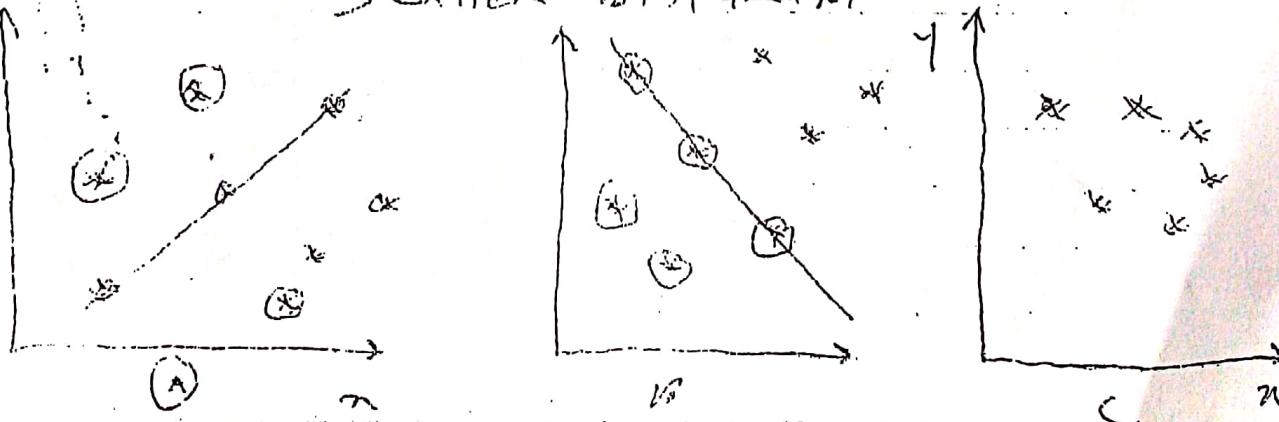
If $r = 0.5$ — 0.59, a moderate or average correlation.

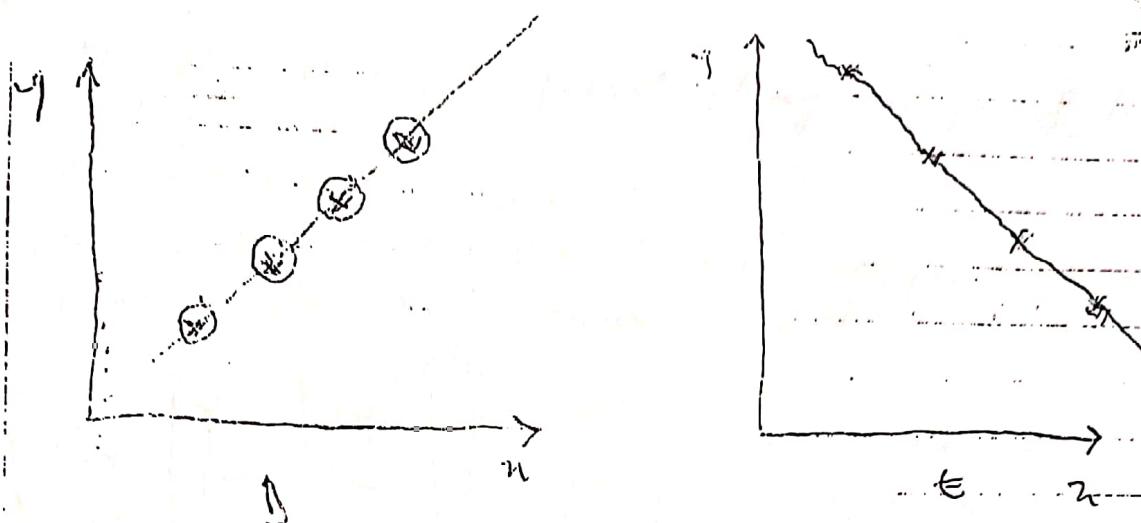
If $r = 0.6$ — 0.99, a strong positive correlation.

If $r = 0$: no correlation,

If $r = 1$, there exists a perfect positive correlation
and vice versa is negative correlation.

SCATTER DIAGRAM





A = positive correlation.

B = Negative Correlation.

C = No Correlation.

D = Perfect Positive.

E = Perfect Negative.

Spearman's Rank Correlation Coefficient is a non-parametric equivalent of the Pearson's product moment correlation coefficient. It's used for data cannot be quantified, the data are rank's first compilation.

The Spearman's Rank formulae is given as

$$\rho_{\text{rank}} = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

Where, $d_i = \text{Rank}(X_i) - \text{Rank}(Y_i)$

$$d_i = U_i - V_i$$

$$\therefore U_i = M_i \quad V_i = N_i \Rightarrow d_i = U_i - V_i = d_i = d_i^2$$

Example 2:

X	15	12	18	21	16
Y	12	14	16	18	18



(i) Plot a scatter diagram for information available.

(ii) Regress y on x .

(iii) Compute the correlation coefficient of x and y and interpret your result.

Solution

X	Y	XY	X^2	Y^2
15	12	180	225	144
12	14	168	144	196
18	16	288	324	256
21	18	378	441	324
16	18	288	256	324

$$\Sigma n = 82, \Sigma y = 78$$

$$\Sigma xy = 1302$$

$$\Sigma x^2 = 1390$$

$$\Sigma y^2 = 1244$$

$$(iv) y = b_0 + b_1 x$$

$$b_1 = \frac{N \Sigma xy - (\Sigma x)(\Sigma y)}{N \Sigma x^2 - (\Sigma x)^2}$$

$$= \frac{5(1302) - 82(78)}{5(1390) - (82)^2}$$

$$b_1 = 0.504$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$\bar{y} = \frac{\Sigma y}{N} = \frac{78}{5} = 15.6$$

$$\bar{x} = \frac{\Sigma x}{N} = \frac{82}{5} = 16.4$$

$$b_0 = 15.6 - (0.504)(16.4)$$

$$b_0 = 7.32$$

$$y = 7.32 + 0.504x$$

$$1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

$$(v) r_{xy} = \frac{5(1302) - (82)(78)}{\sqrt{5(1390) - (82)^2} \cdot \sqrt{5(1244) - (78)^2}}$$

P/S

Date: 11

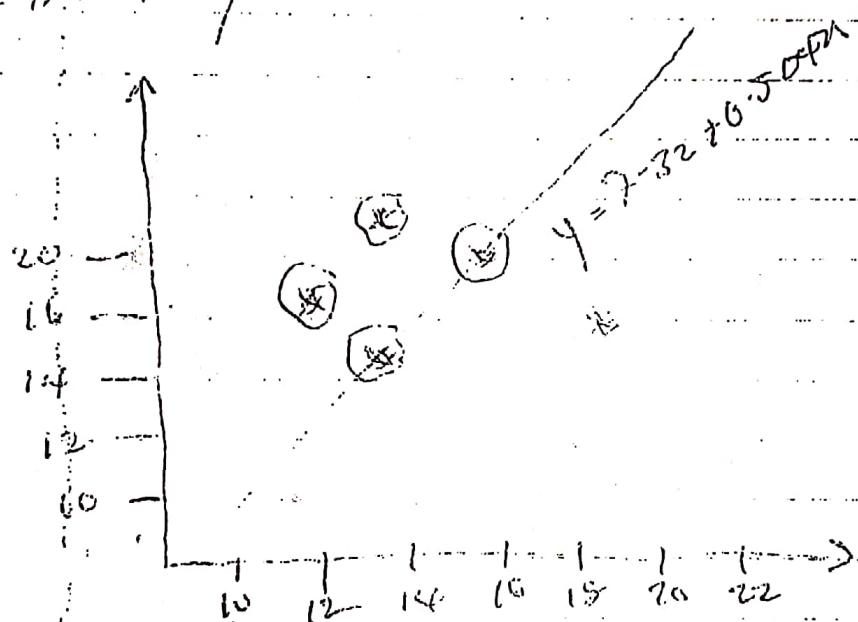
1/20....

$$= \frac{6516 - 6396}{(6.950 - 6.724)(6.220 - 6.084)}$$

$$V_{xy} = \frac{114}{\sqrt{(2.26)(0.36)}} = \frac{114}{\sqrt{303.76}}$$

$$r_{xy} = \frac{114}{14.871} = 0.6371$$

There exist a strong positive correlation between y_1 and y .



Exercise.

- 1) The data below shows the number of road accidents and the resulting deaths in Nigeria between 1990 and 1996.

Year	1990	1991	1992	1993	1994	1995	1996
Number of accidents	17	18	23	25	27	33	54
Deaths	3	4	4	5	5	6	8
Year	1997	1998	1999	2000	2001	2002	2003
Number of accidents	36	36	49	49	54	64	74
Deaths	10	10	14	14	15	16	18

Exponential Smoothing is a rule of thumb technique for fitting time series data using the ~~exponential~~ Window.

Exponential Smoothing is one of many window functions applied to smooth data in signal processing as low pass filters to remove high frequency noise.

- (1) Draw the Scatterplot diagram
 (2) fitted the regression equation line of y on x
 (3) Predict the number of deaths in 1997 if 27 road accident occurred.

(4) The following is a set of data of the federal government direct expenditure on agriculture and health sector in the period of 2001 to 2009.

Period	2001	2002	2003	2004	2005	2006	2007	08	09
Agric.	4.4	4.5	12.4	13.4	24.6	38.8	18.5	42.0	20.5
Cat	12.0	14.1	19.9	18.3	29.1	62.1	83.5	85	95.5

- (5) Estimate the best line of fit between y and X . (iv) What is the degree of association between X and y using simple and Spearman's Correlation coefficient.

TIME SERIES

Is defined as the collection of observation repeated over a long period of time.

Time series is study for the interest, for the future occurrence it also forecast future level economic activities by studying the behavior of data in the past.

COMPONENT OF TIME SERIES

time series is noted for its variation from time to time, this variation can be attributed to climatic, social, economics or accidental factors these factors can be describe as the future of four component of time series described below

- (1) The Trend: It is a part which time series graph appears in column over a long period of time it's upward and downward movement in time series.
- (2) Seasonal Variation: This is brought by climate social factors, Seasonal variation appears with some regularities in the same season of a year, its graph reflects upward and downward movement of time series graph.
- (3) Cyclical Variation: This time series exhibit movement in a fixed period due to some principal cause.
- (4) Irregular Variation: After trend and cyclical variation have been removed from the set of data we are left with series of residual which may or not random they are sporadic or unpredictable like war, fire disaster, flood and election.

MATHEMATICAL REPRESENTATION OF TIME SERIES

The following notations are needed for mathematical combination of time series.

- (a) The time series is denoted by Capital Y.
 - (b) The trend component is denoted by Capital T.
 - (c) The Cyclical Variation is denoted by C.
 - (d) The Seasonal Variation is denoted by Capital S.
 - (e) The irregular variation is denoted by I.
- There are two mathematical model of time series, the Additive and Multiplicative model.
- (a) Additive model is defined by $Y = T + S + C + I$
 - (b) Multiplicative models are:
- $$Y = T \times S \times C \times I = T S C I$$

Exercise.

The data below shows the quantity production of cotton in Nigeria between 2010 and 2013.

Quarters

Year	1	2	3	4
2010	5.8	12.2	16.8	21.1
2011	9.5	22.3	31.0	34.8
2012	20.6	22.5	11.0	40.1
2013	28.7	21.5	11.7	25.5

Compute the trend value using

- (i) the least square method.
- (ii) the moving average.
- (iii) find the seasonal index.

for Seasonal Index - $\frac{\text{Total Value of Quarters}}{\text{Total average}}$

Example

Year	1	2	3	4
2013	13	20	17	25
2014	18	18	16	23
2015	12	22	14	21
2016	11	19	10	24
ifinal	51	79	57	93

$$\text{Quartile average} = \frac{51 + 79 + 57 + 93}{4}$$

$$= \frac{280}{4} = 70$$

$$\text{Seasonal Index for Q1} = \frac{51}{70} = 0.73$$

$$\text{Q1} = 0.73 \times 70 = 51$$



$$\text{Ans} : \bar{x}_3 = \frac{87}{3} = 29.$$

$$\text{Ans} : \bar{x}_4 = \frac{93}{3} = 31.$$

TEST OF HYPOTHESES.

- (1) Statistical hypothesis: is an assertion about the distribution of one or more random variables. If the statistical hypothesis completely specifies the distribution it is called a simple statistical hypothesis. If it doesn't it's called the composite statistical hypothesis. Example.

$$H_0 : \mu = 75$$

H_s

$$H_1 : \mu \neq 75$$

$$H_2 : \mu < 75$$

- (2) A Test of statistical hypothesis is a rule which when experimental sample value have been obtained lead to a decision to accept or reject the hypothesis under consideration.

There are five steps to be taken when one is interested in testing the hypothesis these are:

(i) formulate the hypothesis

(ii) State the level of significance.

$$\text{e.g. } \alpha = 0.05 \text{ or } \alpha = 0.01$$

$$\alpha = 0.1$$

Example: Test at the $\alpha = 0.01$ level significance that the mean of $87, 5$ is less than 6 .

$$H_0 : \mu = 6$$

$$V_s \quad H_1 : \mu < 6$$

0.01 Mean 99% assurance level of the hypothesis you want to take.

(ii) Compute the test statistic (see the formulae).

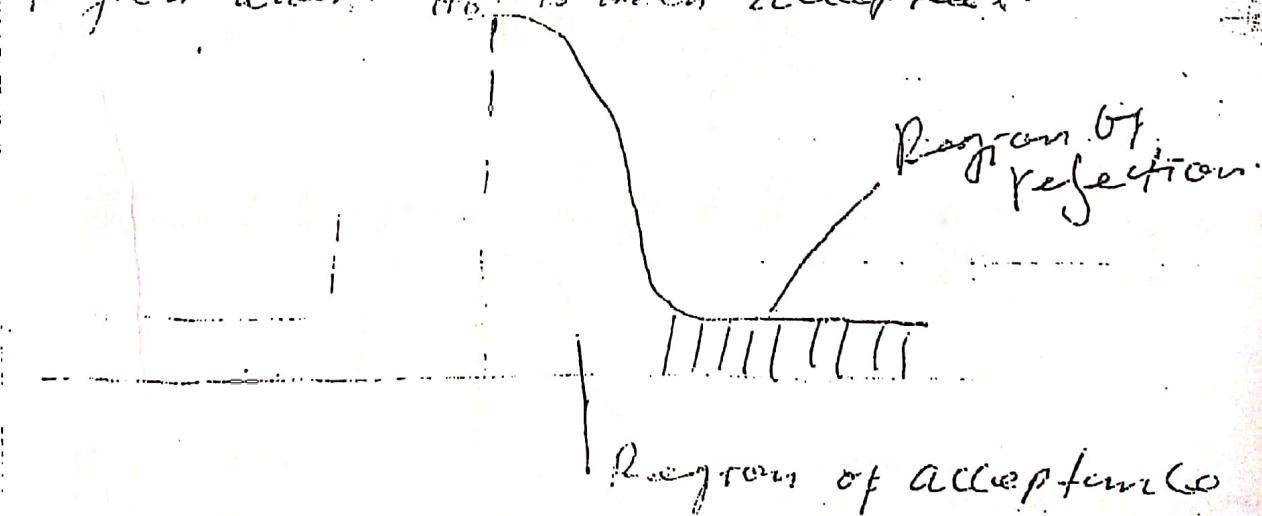
Compare the test statistic value with that in the critical region.

(iii) Arrived at a decision.

5. Level of Significance: it is the quantity of risk of the type I error which we are ready to tolerate by incurring a decision of $\alpha = 0.05$, 0.01 etc.

4. P - Value: this is another approach by which we can find the smallest alpha level at which H_0 is rejected. The criteria for using this is if the p-value is less than α .

(5) CRITICAL REGION: Is a test the area under the probability density curve is divided into two regions regions i.e. the region of rejection and acceptance regions. The region of rejection is the region where H_0 is rejected the region of acceptance is the region where H_0 is accepted.





(6) The T-Test: the critical region is always on the tail of the distribution curve. It may be one-tailed or two-tailed depending upon the alternative hypothesis.

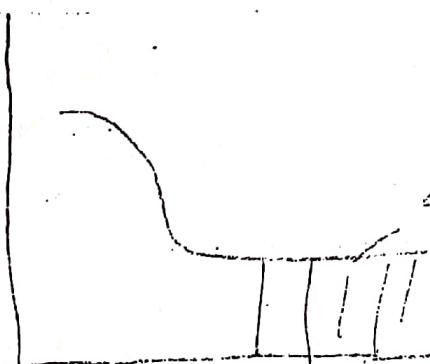
1. One-tailed test (to the right) hypothesis:

$$H_0: \bar{X} = 3$$

$$H_1: \bar{X} > 3$$

$$\text{Test statistic } Z = \frac{\bar{X} - U}{\sigma/\sqrt{n}}$$

Critical region: Reject H_0 if $Z \geq Z_{\alpha}, n-1$



a Rejection Region

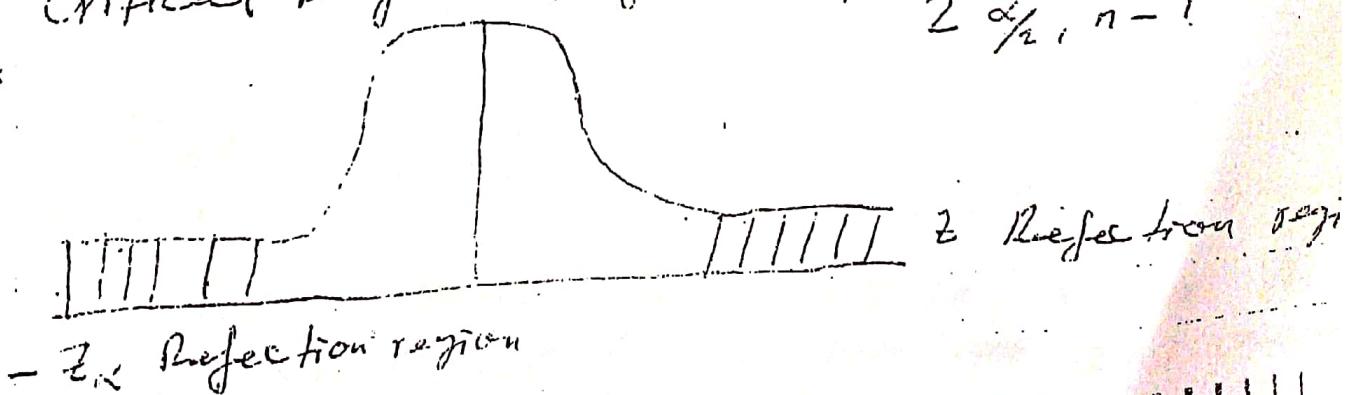
2. One-tailed test (to the left) hypothesis:

$$H_0: \bar{X} = 3$$

$$H_1: \bar{X} < 3$$

$$\text{Test statistic } Z = \frac{\bar{X} - U}{\sigma/\sqrt{n}}$$

Critical Region: Reject H_0 : if $Z \leq -Z_{\alpha}, n-1$ or
 $Z \geq Z_{\alpha/2}, n-1$



- $Z_{\alpha/2}$ Rejection region



Probability of rejecting H_0 when it's false
 This denoted by $1 - \beta$

Example (1) :- The data below represent the scores obtained by 20 students in STA202 and CSC202.

STA202	64	62	52	54	51	61	54	59	57	55
CSC202	66	58	57	68	66	52	63	54	56	58

Assuming the data is normally distributed with means 60 and variance 25. Test at 0.05 level of significance that the mean scores for both exams is less than 60.

Solution

(1) Hypotheses

$$H_0: \mu = 60$$

$$H_1: \mu < 60$$

$$(2) \alpha = 0.05$$

(3) Test statistic

$$Z = \frac{\bar{X} - \mu}{\sqrt{\frac{\sigma^2}{n}}} \quad \text{or} \quad \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$(4) \text{ The means from both } \bar{X} = \frac{\sum x}{n} = \frac{1146}{20}$$

$$\bar{X} = 57.3$$

$$Z = \frac{57.3 - 60}{\sqrt{25}/\sqrt{20}} = \frac{-2.7}{1.12}$$

$$Z_{cal} = -2.42$$

$$Z_{0.05, 19} = -1.73$$

(5) Decision: Since $Z_{cal} = -2.42$ is less than -1.73 we reject the null hypothesis and accept the alternative hypothesis H_1 .

The conclusion reached here is the mean is less than 60.

TESTING DIFFERENCE BETWEEN TWO MEANS (POPULATION VARIANCE ARE UNKNOWN).

formulas:

$$1. \bar{x} = \frac{\sum x_1}{N_1}, \quad 2. \bar{x}_2 = \frac{\sum x_2}{N_2}$$

$$3. S^2 = \frac{1}{N_1 - 1} \sum (\bar{x}_1 - \bar{x})^2 = \frac{N_1 \sum x_1^2 - (\sum x_1)^2}{N_1 - 1}$$

$$4. S^2 = \frac{1}{N_2 - 1} \sum (\bar{x}_2 - \bar{x})^2 = \frac{N_2 \sum x_2^2 - (\sum x_2)^2}{N_2 - 1}$$

$$5. \text{pooled Variance} = \frac{(N_1 - 1) S^2 + (N_2 - 1) S^2}{N_1 + N_2 - 2}$$

$$G. t = \frac{(\bar{x}_1 - \bar{x}_2) - (U_1 - U_2)}{\sqrt{\frac{S^2}{N_1} + \frac{S^2}{N_2}}}$$

Example 2: the information below represent the test scores obtained in 8 for 202 and math 206

9	18	18	3	16	21	13	7	
4	12	5	10	8	11	25	19	17

P. $\alpha = 0.01$ level of significance whether the score if 57 for 202 is equal to the mean of MTH 206.

G.

Solution

1. Hypotheses

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2$$

2. $\alpha = 0.05$

3.

n_1	x_1	x_1^2	n_2	x_2	x_2^2
1	9	81	1	14	196
2	18	324	2	12	144
3	15	225	3	5	25
4	3	9	4	16	100
5	16	256	5	8	64
6	21	441	6	11	121
7	13	169	7	25	625
8	17	289	8	19	361
			9	4	16
			10	17	289

$$n_1 = 8, \sum x_1 = 102, \sum x_1^2 = 1554$$

$$n_2 = 10, \sum x_2 = 125, \sum x_2^2 = 1991$$

$$\bar{x}_1 = \frac{\sum x_1}{n_1} = \frac{102}{8} = 12.75$$

$$\bar{x}_2 = \frac{\sum x_2}{n_2} = \frac{125}{10} = 12.5$$

$$S_1^2 = \frac{n_1 \sum x_1^2 - (\sum x_1)^2}{n_1 - 1} = \frac{8(1554) - (102)^2}{8 - 1}$$

$$S_1^2 = \frac{2028}{7} = 289.7$$

Topic:

$$\begin{aligned} S_p^2 &= \frac{N_2 \sum n_2^2 - (\sum n_2)^2}{N_2 - 1} = \frac{\int_2^2 = 10(1941) - 1125}{10 - 1} \\ &= \frac{3755}{9} = \frac{\int_2^2 = 420.5}{10} \\ S_p^2 &= \frac{(N_1 - 1) s_1^2 + (N_2 - 1) s_2^2}{N_1 + N_2 - 2} = \frac{7(289.7) + 9(420.5)}{16} \end{aligned}$$

$$S_p^2 = 363$$

$$t_{cal} = \frac{12.75 - 2.5}{\sqrt{\frac{363}{8} + \frac{363}{10}}}$$

$$t_{cal} = 0.03$$

4. Decision Rule: Reject H_0 , IF $t_{cal} > t_{\alpha/2, n_1+n_2-2}$

$$t_{0.05, 16} = 2.12$$

5. Since $t_{cal} = 0.03 < t_{tab} = 2.12$ we can't reject the null hypothesis.

6. The conclusion reached there is that the mean score in SF1202 is equal to mean score in MT1206.



MOST POWERFUL TEST

If a simple hypothesis is tested against an alternative and suppose it has a the same level of significance the test will be smaller size for error i.e. the most powerful Test of the two that significant.

A Uniformly most powerful (UMP) is a hypothesis which has the greatest power that is $1 - \beta$ among all the possible test of a given size α .
for example the NEYMAN PEARSON LEMMA
the "likely hood" ratio test etc. are UMP. for testing simple hypotheses.

Best Critical Region: Consider the test of a simple Null hypothesis ($H_0: \theta = \theta_0$) versus a simple alternative $H_1: \theta = \theta_1$. let C be a critical region of size α . the best critical region is that region that has a greatest power among all the critical region of size α .
the Neyman Pearson Lemma gives sufficient conditions for a best critical region of size α .

Neyman Pearson Lemma:
let x_1, x_2, \dots, x_n be a random sample of size n from its distribution with PDF $f(x; \theta)$ where θ_0 and θ_1 are two possible values of θ . let the joint PDF of x_1, x_2, \dots, x_n be denoted by the likely hood function $L(\theta) = L(\theta; x_1, x_2, \dots, x_n) = f(x, \theta)$

$f(x, \theta) = f(x_1, \theta) \cdots f(x_n, \theta)$ if and only if (IFF)
there exist a positive constant k and a critical region C of the sample space such that



TUTORIAL

MOST POWERFUL TEST

If a simple hypothesis is tested against its alternative and suppose it has a the same level of significance. The test will be smaller size for error if the most powerful Test of the two that significant.

A Uniformly most powerful (UMP) is a hypothesis which has the greatest power that is $1 - \beta$ among all the possible test of a given size α . for example the NEYMAN PEARSON LEMMA the likelihood ratio test etc. are UMP. for testing simple hypotheses.

Best Critical Region: Consider the test of a simple Null hypothesis ($H_0: \theta = \theta_0$) versus a simple alternative $H_1: \theta = \theta_1$. Let C be a critical region of size α . The best critical region is that region that has a greatest power among all the critical region of size α .

The Neyman Pearson Lemma gives sufficient conditions for a best critical region of size α .

NEYMAN PEARSON LEMMA:

Let x_1, x_2, \dots, x_n be a random sample of size n from its distribution with PDF $f(x; \theta)$ where θ_0 and θ_1 are two possible values of θ . Let the joint PDF of x_1, x_2, \dots, x_n be denoted by the likelihood function $L(\theta) = L(\theta; x_1, x_2, \dots, x_n) = f(x, \theta)$.

$f(x_1, \theta_0) \cdots f(x_n, \theta_0)$ if and only if (IFF) there exist a positive constant k and critical region C of the sample space such that

- (i) $P(x_1, x_2, \dots, x_n) \in C; \Omega_0) = \alpha$
- (ii) $\frac{L(\Omega_0)}{L(\Omega_1)} \leq K$ for $(x_1, x_2, \dots, x_n) \in C$
- (iii) $\frac{L(\Omega_1)}{L(\Omega_0)} \geq k$ for $(x_1, x_2, \dots, x_n) \in C'$

Here C is a best critical region of size α .
 i.e., testing the simple hypothesis ($H_0: \theta = \theta_0$) against an alternative simple hypothesis ($H_1: \theta = \theta_1$).
 Note that a best critical region can be determined by using the Neyman Pearson Lemma if and only if

$$\frac{L(\Omega_0)}{L(\Omega_1)} \leq K \text{ if } \bar{x} \in C$$

$$\text{or} \\ \frac{L(\Omega_1)}{L(\Omega_0)} \geq k \text{ if } \bar{x} \in C$$

Example. Let x_1, x_2, \dots, x_n be a random sample from a normal distribution with $N(u, \sigma^2)$. Find the best critical region for testing the $H_0: u = 55$.

Solution.

$$H_0: u = 55$$

$$f(x; \theta) = (2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2} \sum (x_i - u)^2}$$

$$L(x; \theta) = (2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2} \sum (x_i - u)^2}$$

$$\frac{L_0}{L_1} \leq K$$



$$\frac{L(50)}{L(55)} = \frac{(72\pi)^{-\frac{n}{2}}}{(72\pi)^{-\frac{n+1}{2}}} = \exp\left(-\frac{1}{72}\right)$$

$$\exp\left[-\frac{1}{72} \left\{ (x_1 - 50) + (x_2 - 55) \right\}\right]$$

$$\frac{L(50)}{L(55)} = \exp\left(-\frac{1}{72}\right) \cdot \frac{\mathbb{E}(x_1^2 - 100x_1 + 2500)}{\mathbb{E}(x_1^2 - 110x_1 + 3025)}$$

$$= \exp\left(-\frac{1}{72}\right) \frac{\mathbb{E}x_1^2 - 100\mathbb{E}x_1 + 2500}{\mathbb{E}x_1^2 - 110\mathbb{E}x_1 + 3025}$$

$$= \exp\left[-\frac{1}{72}\right] \left[\mathbb{E}x_1^2 - 100\mathbb{E}x_1 + 2500 - \mathbb{E}x_1^2 + 110\mathbb{E}x_1 - 3025 \right] \leq k.$$

$$= \exp\left[-\frac{1}{72}\right] [10\mathbb{E}x_1 - 5525] \leq k$$

If we take the natural log of each member of inequality we find out that $-\frac{1}{72} (20\mathbb{E}x_1 - 5525) \leq \ln k$

$$\leq \ln k$$

Now multiplying through by 72

$$(10\mathbb{E}x_1 - 5525) \leq 72 \ln k$$

To make $\mathbb{E}x_1$ the subject we have

$$-10\mathbb{E}x_1 \leq 72 \ln k - 5525$$

$\mathbb{E}x_1 \geq 52.5 - 7.2 \ln k$ by dividing both side by 10

$$\bar{x} \geq 52.5 - \frac{7.2 \ln k}{n}$$

$$\bar{x} \geq c$$

where $c = 52.5 - 7.2 \frac{\ln k}{n}$ is the best critical region according to NPL.

PTB

9/9/2020

Date: / /

1/20....



Topic:

Exercise:

- (1) Use Person Lemma (NPC) to obtain the best critical region for testing $H_0: \theta = 0$.

$H_0: \theta = 0$,
for a normal distribution given as
 $X \sim N(\mu, \sigma^2 I)$

Example: determine whether the following vector are

- (2) Show that $C = \{x_1, x_2, \dots, x_n\} \times L_C$ is a best critical region for testing $H_0: \mu = 80$

$$C = \{x_1, x_2, \dots, x_n\} \times L_C$$

L_C for $N(\mu, 64)$



UNBIASED TEST.

In statistical hypothesis testing, a test is said to be unbiased if the probability of committing type I error is equal to the significance level.

Unbiased estimator: An estimator $\hat{\theta}$ of parameter θ is unbiased if the expected value $\theta = \hat{\theta} = E(\hat{\theta})$ = some distribution.

$$(i) \text{ Bernoulli} = f(x|\theta) = \theta^x (1-\theta)^{1-x} \quad E(X) = \theta, \text{Var} = \theta(1-\theta).$$

$$(ii) \text{ Poisson} = f(x|\theta) = \frac{\theta^x e^{-\theta}}{x!} \quad E(X) = \theta, \text{Var} = \theta$$

$$(iii) \text{ Binomial} \quad B(X; n, \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \quad E(X) = n\theta, \text{Var}(X) = n\theta(1-\theta)$$

(iv) Exponential

$$f(x|\theta) = \frac{1}{\theta} e^{-x/\theta}$$

$$E(X) = \theta \quad \text{Var}(X) = \theta^2$$

(θ) plan

(v) Normal

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

Example: Let X_1, X_2, \dots, X_n be a random sample from an exponential distribution with parameter θ .

- (i) Show that $\hat{\theta}$ is an unbiased estimator of θ i.e. $E(\hat{\theta}) = \theta$.
- (ii) Show that the variance $\text{Var}(\hat{\theta}) = \frac{\theta^2}{n}$.

Solution:

$$f(x|\theta) = \begin{cases} \frac{1}{\theta} e^{-x/\theta} & 0 \leq x \leq \infty \\ 0, & \text{elsewhere} \end{cases}$$

(i) Mean of exponential distribution.

$$\begin{aligned} E(X) &= \hat{\theta} \\ \therefore \hat{\theta} &= \bar{x} \\ E(\hat{\theta}) &= E(\bar{x}) \\ E(\bar{x}) &= E\left(\frac{\sum x}{n}\right) \\ &= \frac{1}{n} \sum E(x) \end{aligned}$$

$$\begin{aligned} \text{Since } E(x) &= \theta \\ &= \frac{1}{n} \sum (\theta) \\ &= \frac{1}{n} \cdot n \cdot \theta \end{aligned}$$

$$E(\hat{\theta}) = \theta.$$

(ii) Variance of exponential distribution

$$\begin{aligned} V(\bar{x}) &= V\left(\frac{\sum x}{n}\right) \\ &= \frac{1}{n^2} \sum V(x) \end{aligned}$$



$$\geq \frac{1}{n^2} \sum (\text{Var}(x_i))$$

$$\text{Var}(n) = \theta^2$$

$$= \frac{1}{n^2} \sum (\theta^2)$$

$$\Rightarrow \frac{1}{n^2} \cdot n \theta^2 = \frac{\theta^2}{n}$$

$$\therefore \text{Var}(\bar{x}) = \frac{\theta^2}{n}$$

Exercise:

- ① Let x_1, x_2, \dots, x_n be a random sample from binomial distribution with parameter n and (θ) , $B(x; n, \theta)$. Find the ~~MLE~~ N.M. and show that \bar{x} is an unbiased estimator of θ .

Likelihood Ratio Test for Univariate Distribution: Let $L(\theta)$ equal to $L(x_1, x_2, \dots, x_n; \theta)$ be a likelihood function for the random variables x_1, x_2, \dots, x_n . If $\hat{\theta}$ is the value of θ which maximizing likelihood function then $\hat{\theta}$ is called the maximum likelihood estimator (MLE) of θ for the sample x_1, x_2, \dots, x_n .

Steps/Conclusion

- ② Many likelihood function satisfied regularity condition. The L.M.E is the solution to the equation

$$\frac{\partial L(\theta)}{\partial \theta} = 0$$

Dated 11

J/20...



② $L(\theta)$ and natural log $\log_e L(\theta)$ has been found to have the maximum at the same value θ and it's sometimes easier to find the maximum of natural log $\log_e L(\theta)$.

Note if x_1, x_2, \dots, x_n is a random sample, the "likelihood" function is the joint PAF $f(x_1, x_2, \dots, x_n | \theta)$ which is given as

$$L(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta).$$

Maximum Likelihood Estimator (MLE): An estimator $\hat{\theta}$ is defined as a maximum likelihood estimator (MLE) if

$$L(x_1, x_2, \dots, x_n | \theta) = \text{Max}_{\theta} \{ L(x_1, x_2, \dots, x_n | \theta) \}$$

Example 1: Given a random sample x_1, x_2, \dots, x_n of a binomial θ random sample. Find the MLE of θ .

Solution:

$$f(x; \theta) = \theta^x (1-\theta)^{1-x}$$

$$L(x; \theta) = \prod_{i=1}^n \{ \theta^{x_i} (1-\theta)^{1-x_i} \}$$

$$= \prod_{i=1}^n \ln \{ (1-\theta)^{n-x_i} \}$$

$$\text{where } t = \sum_{i=1}^n (1-\theta)^{n-x_i}$$

$$= \theta^t (1-\theta)^{n-t}$$

$$\log_e L(\theta, \theta) = t \log \theta + n - t \log(1-\theta)$$

1
x)



Set
 $\frac{\partial \log L(\theta; \mathbf{x})}{\partial \theta} \Rightarrow 0$

$$\therefore \frac{t}{\theta} - \frac{n-t}{1-\theta} = 0$$

$$\frac{t}{\theta} = \frac{n-t}{1-\theta}$$

$$t(1-\theta) = \theta(n-t)$$

$$t - t\theta = \theta n - \theta^2 t$$

$$t - \theta n = t\theta = \theta^2 t$$

$$\hat{\theta} = \frac{t}{n}, \quad \theta = \frac{tn}{n} = \bar{x}$$

Example 2: Find the maximum likelihood estimate of Poisson distribution with parameter θ

Solution:

$$f(x_i; \theta) = \frac{\theta^n e^{-\theta}}{x_i!}$$

$$L(\mathbf{x}; \theta) = \prod_{i=1}^n \frac{(\theta^n e^{-\theta})}{x_i!}$$

$$= \frac{\theta^{n\bar{x}} e^{-n\theta}}{\prod_{i=1}^n x_i!}$$

where $\prod_{i=1}^n x_i = \text{error}$ $\bar{x} = \bar{e}$

$$= \theta^{\bar{x}} e^{-n\theta}$$



$$\log_e L(x; \theta) = t \log \theta - n\theta.$$

$$\frac{\partial \log_e L(x; \theta)}{\partial \theta} = 0$$

$$\frac{\ell}{\theta} - n = 0$$

$$\frac{\ell}{\theta} = n, \quad t = n\theta, \quad \theta = \frac{t}{n}$$

$$\hat{\theta} = \frac{tn}{n}, \quad \hat{\sigma}^2 = \bar{x}^2$$

Exercise

- (i) Given that a random sample x_1, x_2, \dots, x_n of a normal distribution $X \sim N(\mu, \sigma^2)$
- The MLE of μ and σ^2
 - Show that \bar{x} is an unbiased estimator of μ .
 - Show that $\hat{\sigma}^2$ is not an unbiased estimator of σ^2 .



Topic: Exercise,

Bernoulli

$$(x_1, \dots, x_n) = \frac{\theta^n (1-\theta)^{n-x}}{A_n}$$

$$= \frac{\theta^n (1-\theta)^{n-x}}{A_n}$$

$$= \log \theta^n + \log (1-\theta)^{n-x} =$$

$$\text{Neffektiv} = \text{En log } \theta + \frac{(n-x)}{\theta} \log (1-\theta)$$

$$\text{Effektiv} = \text{EnL} - \frac{x}{1-\theta} = 0$$

$$\text{En} = \frac{n - x}{1-\theta}$$

$$\text{En} - \text{Ex} = \text{En} - \theta n$$

$$\theta n = \text{Ex}$$

$$x = \frac{\text{Ex}}{n}$$

$$\text{Binomial} = \theta^n (1-\theta)^{n-x}$$

$$P(x) = L = \frac{\text{En}^x}{A_n} (1-\theta)^{n-x}$$

$$= \text{En} \log \theta + (n - \text{En}) \log (1-\theta)$$

Poisson Distribution.

$$p(x/\theta) = \frac{\theta^x e^{-\theta}}{x!}$$

Normal distribution

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



Youden Squares:

These are incomplete latin square designs in which the number of columns does not equal the number of rows and treatments.

Consider the design shown in the table below for Youden Squares for five treatments (A, B, C, D, E).

Row

is nothing it's not necessary to that every Latin Square with more than one column

the below model of Youden-Square is

$$Y_{ijk} = \bar{M} + \alpha_i + T_j + B_h + \epsilon_{ijk}$$

where

\bar{M} is the overall mean

α_i is the i th block effects to j th treatment effect, B_h is the h th position error term or error, T_j is the j th treatment effect and ϵ_{ijk} is the $N(0, \sigma^2)$ error since position occurs exactly once position are orthogonal to treatment and blocks.

example:

Day (Block)	1	2	3	4	Σ_i
1	$A=3$	$B=1$			
2	$B=0$	$C=0$			
3	$D=1$	$D=0$			
4	$D=1$	$E=0$			
5	$E=5$				