

STATISTICS

18/04/2021

→ eg - Budget

Data: A number in either ascending or descending order.

Class interval: lower value, upper value

Def: Statistics is defined as a branch of science, dealing with collection of data, organising, summarising, presenting and analysing data and drawing a valid conclusion.

And thereafter making reasonable decisions on the basis of such analysis.

eg → Analysing the data we can conclude that whether rain will occur or not.

Discrete data:

eg There is a break

Continuous data:

eg

Continuously	0-10	5
	10-20	8
	20-30	4

CURVE FITTING:

The general problem of finding equations of approximating curves that closely fit the given data is called curve fitting.

In other words, curve-fitting is the representation of relationship b/w the ^{two} variables by means of an algebraic expression.

- This type of suggestion first given by scatter.
- And the diagram is known as scattered diagram
- calculation of scattered diagram by Li-square method.

* The curve passes through maximum no. of points from scattered diagram, more than one curve may be seen to be appearing to the given set of data and this may largely be due to human judgement.

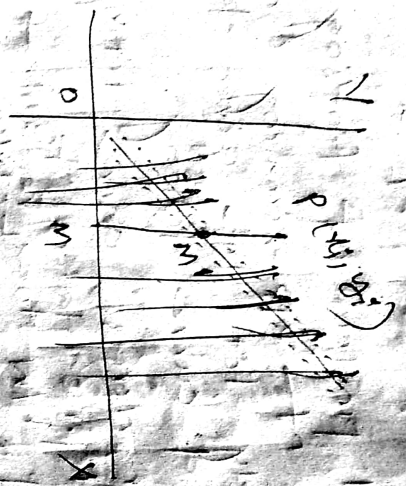
* To achieve this method of Li-square is used

* let $P(x_i, y_i)$ be one of the points, of the dotted diagram

* from P draw \perp PM on the X axis

* let it meet the curve at M'

* The distance $PM' = PM - MM'$



* p_n is called the error of estimate or residual or deviation.

* This may be positive or zero or negative depending upon whether y lies above, on or below

axis the curve.

* Since residuals corresponding to $(n-1)$ points may be obtained.

* Let the residuals be denoted by $e_1, e_2, e_3, \dots, e_{n-1}$

* Sum of the squares of these residuals, i.e.

$$e_1^2 + e_2^2 + e_3^2 + \dots + e_{n-1}^2$$

$$\sum e_i^2$$

$$\sum e_i^2$$

$$\sum e_i^2$$

* If $\sum e_i^2$ is small then fit is good. If $\sum e_i^2$ is large then fit is poor.

* This method is known as least square method.

The line of best fit is

Let $y = a + bx$ be the line of best fit.

Let (x_i, y_i) be n points through which it passes.

If these lines will be lines of best fit then

$$\sum_{i=1}^n e_i^2 = E = \sum_{i=1}^n (y_i - a - bx_i)^2$$

* E to be maxima or minima:

By using principle of n square:

Partial derivative of E w.r.t a and b must vary separately

$$\frac{\partial E}{\partial a} = -2 \sum (y_i - a - bx_i) = 0$$

$$\frac{\partial E}{\partial b} = -2 \sum (x_i)(y_i - a - bx_i) = 0$$

$$\sum y_i - a - bx_i = 0$$

$$\sum x_i (y_i - a - bx_i) = 0$$

$$\Rightarrow \sum y_i = na + b \sum x_i$$

$$\sum x_i y_i = a \sum x_i + b \sum x_i^2$$

$$a) y = ax + b$$

$$\sum y = a \sum x + nb$$

$$\sum xy = a \sum x^2 + b \sum x$$

$$b) y = a + bx + cx^2$$

$$E = \sum (y - a - bx - cx^2)^2$$

$$\frac{\partial E}{\partial a} = -2 \sum (y - a - bx - cx^2) = 0$$

$$\frac{\partial E}{\partial b} = -2 \sum x (y - a - bx - cx^2) = 0$$

$$\frac{\partial E}{\partial c} = -2 \sum x^2 (y - a - bx - cx^2) = 0$$

$$\sum y = na + b \sum x + c \sum x^2$$

$$\sum xy = a \sum x + b \sum x^2 + c \sum x^3$$

$$\sum x^2 y = a \sum x^2 + b \sum x^3 + c \sum x^4$$

$$a) y = a + bx + c$$

$$\sum y = a \sum 1 + b \sum x + c \sum 1$$

$$\sum xy = a \sum x + b \sum x^2 + c \sum x$$

$$\sum x^2 y = a \sum x^2 + b \sum x^3 + c \sum x^2$$

Let a straight line $y = a + bx$ by method of least square from the following data points.

x	0	1	3	6	8
y	1	3	2	5	4

Q.1 $y = a + bx + cx^2$

Q.2 $y = a + bx$

$y = a + bx$

$\sum y = na + b \sum x$

$\sum xy = a \sum x + b \sum x^2$

x	y	x^2	xy
0	1	0	0
1	3	1	3
3	2	9	6
6	5	36	30
8	4	64	32

$\sum x = 18$

$\sum x^2 = 110$

$\sum y = 15$

$\sum xy = 71$

$15 = 5a + b \times 18$

$71 = 18a + 110b$

$y = 1.64 + 0.37x$

solving we get
 $a = 1.64$
 $b = 0.376$

x	y	x^2	x^3	x^2y	xy
2	10	4	8	20	20
0	1	0	0	0	0
1	3	1	1	3	3
3	2	9	27	18	6
6	5	36	216	180	30
8	4	64	512	256	32
18	15	324	756	4572	270

$$y = a + bx + cx^2$$

$$\sum y = na + b \sum x + c \sum x^2$$

$$\sum xy = a \sum x + b \sum x^2 + c \sum x^3$$

$$\sum x^2y = a \sum x^2 + b \sum x^3 + c \sum x^4$$

$$15 = 5a + 18b + c \times 110$$

$$71 = a \times 18 + b \times 110 + c \times 756$$

$$457 = 110a + 756b + 4572c$$

$$a = 1.34$$

$$b = 0.34 \times 10^{-3}$$

$$c = 44.97 \times 10^{-3}$$

$$a = 1.34$$

$$b = 0.34 \times 10^{-3}$$

$$c = -0.0041 \times 10^{-3}$$

$$y = 1.34 + 0.34x - 0.04x^2$$

FITTING OF OTHER CURVES

$$y = ax^b$$

$$\log_{10} y = \log_{10} a + b \log_{10} x$$

$$Y = A + BX$$

$$Y = \log_{10} y$$

$$A = \log_{10} a$$

$$B = \log_{10} x$$

$$y = ae^{bx}$$

$$\log_{10} y = \log_{10} a + b \log_{10} e$$

$$\log_{10} y =$$

$$Y = A + BX$$

$$B = b \log_{10} e$$

Combined gas equation

$$PV^r = \text{const.}$$

$$r y^n = b$$

$$\log_{10} x + r \log_{10} y = \log_{10} \frac{b}{r}$$

$$\log y = \frac{1}{a} \log b - \frac{1}{a} \log_{10} x$$

$$x = 10 \log_{10} y$$

$$y = 10 \log_{10} x$$

$$A = 1/9 \log_{10} 10$$

$$B = -\frac{1}{a} \log_{10} 10$$

$$a \log_{10} y = x$$

Find the curve of best fit of the type $y = a + bx$

data by method of least square

Data given:

x	1	2	3	4	5	6	7	8	9	10
y	1.1	1.5	2.0	2.5	3.0	3.5	4.0	4.5	5.0	5.5

$$y = a + bx$$

$$y = \log_{10} x$$

$$A = \log_{10} a$$

$$B = b \log_{10} e$$

x	y	$\log_{10} x$	y/x	$1/x$
1	10	0	10	1
2	15	0.3010	7.5	0.5
3	12	0.4771	4	0.3333
4	9	0.6021	2.25	0.25
5	8	0.6990	1.6	0.2
6	7	0.7782	1.1667	0.1667
7	6	0.8451	0.8571	0.1429
8	5	0.9031	0.625	0.125
9	4	0.9542	0.4444	0.1111
10	3	1.0	0.3	0.1
Σx	55			
Σy	73			
Σx^2		130		
Σy^2		556		

$$y = a + b \log_{10} x$$

$$\Sigma y = n a + b \Sigma \log_{10} x$$

$$\Sigma xy = a \Sigma x + b \Sigma x^2$$

$$573 = 5a + 34b$$

$$556 = 34a + 300b$$

$$a = 8.70$$

$$b = 0.967$$

$$y = 8.70 + 0.967 \log_{10} x$$

$$y = 8.70 + 0.967 \log_{10} x$$

Co-relation and Regression:

→ Co-relation is the co-variation of 2 independent magnitudes.

If 2 variables x and y are related in such a way that increase or decrease in one of them correspond to increase or decrease in the other.

→ In other words we say that the variables are positively co-related, also if increases or decrease in one of them correspond to decrease or increase in the other, then the

variables are said to be negatively co-related.

→ The numerical of co-relation b/w 2 variables x and y is known as

Pearson's co-efficient of co-relation, usually denoted by r .

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

$$\sigma_x = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

$$= \frac{\sum x^2}{n}$$

$$X = x - \bar{x}$$

$$\sigma_y = \sqrt{\frac{\sum (y - \bar{y})^2}{n}}$$

$$Y = y - \bar{y}$$

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}}$$

→ The coefficient of correlation numerically does not exceed unity.

$$-1 \leq r \leq 1$$

→ If $r = 1$, we say that x and y are perfectly co-related.

→ If $r = 0$, we say that x and y are non-co-related.

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} = \frac{\sum xy}{\sqrt{n \cdot \sigma_x^2} \sqrt{n \cdot \sigma_y^2}} = \frac{\sum xy}{n \sigma_x \sigma_y}$$

Regression

→ It is an estimation of one independent variable in terms of other.

→ If x and y are co-related the best fitting straight line in the

least square sense gives a reasonable good relation b/w x and y .

→ The best fitting straight line of the form $y = ax + b$

→ (as being the independent variables) is called the regression

of y on x :
 → again if $x = 10y + 5$ (any be the independent variable) is called regression of x on y .

→ Derivation of eqn of the regression line:

Regression line of y on x :

$$(y - \bar{y}) = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

Regression of x on y :

$$(x - \bar{x}) = r \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

Problem 1 Compute co-efficient of correlation and the eqn of the lines of regression from the following data:

x :	10	2	3	4	5	6	7
y :	9	8	10	12	11	13	14

$$\sigma_x^2 + \sigma_y^2 = \sigma_{xy}^2$$

$$2\sigma_x \sigma_y$$

We shall prepare the following columns:

x	y	$\sum x-y$	x^2	y^2	xy
1	3	-2	1	9	3
2	3	-1	4	9	6
3	5	-2	9	25	15
4	5	-1	16	25	20
5	5	0	25	25	25
6	5	1	36	25	30
7	5	2	49	25	35
8	5	3	64	25	40
9	5	4	81	25	45
10	5	5	100	25	50
11	5	6	121	25	55
12	5	7	144	25	60
13	5	8	169	25	65
14	5	9	196	25	70
15	5	10	225	25	75
16	5	11	256	25	80
17	5	12	289	25	85
18	5	13	324	25	90
19	5	14	361	25	95
20	5	15	400	25	100
21	5	16	441	25	105
22	5	17	484	25	110
23	5	18	529	25	115
24	5	19	576	25	120
25	5	20	625	25	125
26	5	21	676	25	130
27	5	22	729	25	135
28	5	23	784	25	140
29	5	24	841	25	145
30	5	25	900	25	150
31	5	26	961	25	155
32	5	27	1024	25	160
33	5	28	1089	25	165
34	5	29	1156	25	170
35	5	30	1225	25	175
36	5	31	1296	25	180
37	5	32	1369	25	185
38	5	33	1444	25	190
39	5	34	1521	25	195
40	5	35	1600	25	200
41	5	36	1681	25	205
42	5	37	1764	25	210
43	5	38	1849	25	215
44	5	39	1936	25	220
45	5	40	2025	25	225
46	5	41	2116	25	230
47	5	42	2209	25	235
48	5	43	2304	25	240
49	5	44	2401	25	245
50	5	45	2500	25	250
51	5	46	2601	25	255
52	5	47	2704	25	260
53	5	48	2809	25	265
54	5	49	2916	25	270
55	5	50	3025	25	275
56	5	51	3136	25	280
57	5	52	3249	25	285
58	5	53	3364	25	290
59	5	54	3481	25	295
60	5	55	3600	25	300
61	5	56	3721	25	305
62	5	57	3844	25	310
63	5	58	3969	25	315
64	5	59	4096	25	320
65	5	60	4225	25	325
66	5	61	4356	25	330
67	5	62	4489	25	335
68	5	63	4624	25	340
69	5	64	4761	25	345
70	5	65	4900	25	350
71	5	66	5041	25	355
72	5	67	5184	25	360
73	5	68	5329	25	365
74	5	69	5476	25	370
75	5	70	5625	25	375
76	5	71	5776	25	380
77	5	72	5929	25	385
78	5	73	6084	25	390
79	5	74	6241	25	395
80	5	75	6400	25	400
81	5	76	6561	25	405
82	5	77	6724	25	410
83	5	78	6889	25	415
84	5	79	7056	25	420
85	5	80	7225	25	425
86	5	81	7396	25	430
87	5	82	7569	25	435
88	5	83	7744	25	440
89	5	84	7921	25	445
90	5	85	8100	25	450
91	5	86	8281	25	455
92	5	87	8464	25	460
93	5	88	8649	25	465
94	5	89	8836	25	470
95	5	90	9025	25	475
96	5	91	9216	25	480
97	5	92	9409	25	485
98	5	93	9604	25	490
99	5	94	9801	25	495
100	5	95	10000	25	500

$$\bar{x} = \frac{\sum x}{n} = \frac{28}{7} = 4$$

$$\bar{y} = \frac{\sum y}{n} = \frac{14}{7} = 2$$

$$\sum (x - \bar{x}) = 0$$

$$\sum x^2 = \frac{\sum x^2}{n} = \frac{140}{7} = 20$$

$$\sum y^2 = \frac{\sum y^2}{n} = \frac{14}{7} = 2$$

$$\sum xy = \frac{\sum xy}{n} = \frac{84}{7} = 12$$

$$\sum x^2 = \frac{\sum x^2}{n} = \frac{140}{7} = 20$$

$$r = \frac{4+4-4/2}{2 \times 2 \times 2}$$

$$r = 0.928$$

$$= 0.93$$

Regression of y on x :

$$(y - \bar{y}) = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$= (y - 11) = 0.93 (x - 4)$$

$$= y - 11 = 0.93x - 3.72$$

$$= 0.93x - y + 7.28$$

Regression of x on y :

$$(x - \bar{x}) = r \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$= (x - 4) = 0.93 (y - 11)$$

$$= x - 4 = 0.93y - 10.23$$

$$= x - 0.93y + 6.23$$

Given :
 $8x - 1$
 $40x -$
 are the
 find
 Co-ord
 find
 Ans: Suppose
 the
 $8x$
 $40x$
 40
 $8x$
 40

Given:

$$8x - 10y + 66 = 0$$

$$40x - 18y = 214$$

are the 2 regression lines, find
the means of x and y and the
co-relation coefficient.

$$\text{Find } \bar{y} \text{ of } 8x = 0$$

Ans. Suppose 2 regression lines passing through
the means, i.e. they must satisfy the eqn.

$$8\bar{x} - 10\bar{y} + 66 = 0$$

$$40\bar{x} - 18\bar{y} = 214$$

$$\bar{x} = 13, \bar{y} = 17$$

$$8x - 10y + 66 = 0$$

$$10y = 8x + 66$$

$$y = 0.8x + 6.6 \quad \text{on } y$$

$$40x - 18y = 214$$

$$x = \frac{214}{40} + \frac{18y}{40} \quad \text{on } x$$

$$= 0.45y + 5.35$$

$$r \frac{\bar{y}}{\bar{y}} = 0.8 \times \frac{\bar{y}}{\bar{y}} = 0.45$$

$$r = 0.6$$

← as r is always ± 1 .

Whole those of sample are denoted by \bar{x} and s for mean & variance

Standard error:

The standard deviation of the sampling distribution of a statistic is known as the standard error.

→ It plays an important role in the theory of large sample and it forms the basis of testing of hypothesis.

→ for any statistic for large sample

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

is normally distributed with mean μ and variance $\frac{\sigma^2}{n}$

→ for large sample the theory of

→ Sample size

→ Population variance

→ Sample variance

→ Population proportion

→ Sample proportion

→ Sample size

→ Sample size

→ Sample size

Sl. No	Statistic	S.O.E
1	\bar{x}	$\frac{\sigma}{\sqrt{n}}$
2	s	$\frac{\sigma}{\sqrt{n}}$
3	Difference of 2 sample means $\bar{x}_1 - \bar{x}_2$	$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$
4	Difference of 2 sample standard deviations $s_1 - s_2$	$\sqrt{\frac{\sigma_1^2}{2n_1} + \frac{\sigma_2^2}{2n_2}}$
5	Difference of 2 sample proportions $p_1 - p_2$	$\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$
6	Observed sample proportions p	$\sqrt{\frac{pq}{n}}$

Test of significance:

Test of significance occurs in 2 conditions:

1. The deviation b/w the observed sample statistics and the hypothetical parameter.
2. The deviation b/w 2 sample statistics is significant due to the fluctuation of sampling.

Testing of hypothesis:

Hypothesis is a definite statement as the population parameter for applying test of significance. It is called null hypothesis.

denoted by H_0

Alternative hypothesis:

Any hypothesis complementary to null hypothesis is called alternative hypothesis, denoted by H_1

We want to test null hypothesis H_0 with mean μ_0 under the specified condition

$$H_0: \mu = \mu_0$$

Alternative hypothesis will be as follows

$$H_1: \mu \neq \mu_0 \quad \left\{ \begin{array}{l} \mu > \mu_0 \\ \mu < \mu_0 \end{array} \right.$$

Two tailed alternative hypothesis

$$\mu \neq \mu_0$$

Right tailed alternative hypothesis or single tailed

$$\mu > \mu_0$$

Left tailed alternative hypothesis or single tailed

Alternative hypothesis tells us whether the test is two tailed or one tailed test

Critical region

A region corresponding to a statistic T in the sample space S of which amounts to rejection of the null hypothesis H_0 is called as critical region.

→ The region of the sample space which amounts to acceptance of H_0 is called acceptance region.

→ Critical region is also known as region of rejection.

$$P(Z_1 > Z_\alpha) = \alpha$$

$$P(Z > Z_\alpha) = \alpha$$

$$P(Z < -Z_\alpha) = \alpha$$

Level of significance:

	1% (0.01)	5% (0.05)	10% (0.10)
Two-tailed test	$ Z = 2.58$	$ Z = 1.96$	$ Z = 1.64$
Right-tailed test	$Z_\alpha = 2.33$	$Z_\alpha = 1.645$	$Z_\alpha = 1.28$
Left-tailed test	$Z_\alpha = -2.33$	$Z_\alpha = -1.645$	$Z_\alpha = -1.28$

working procedure for testing of statistical hypothesis.

Step 1) Null hypothesis?

Set up H_0

Step 2) Alternative hypothesis?

Set up H_1

So that we can decide whether we should use one-tailed system or two-tailed system.

Step 3) level of significance

Select the appropriate level of significance in advance depending on the reliability of the estimate

Step 4) Test of statistics

Compute the test statistic

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

$$S.E(\bar{x})$$

Under the null hypothesis.

Conclusion

Step 5)

Compare the computed

value of Z with critical values at level of

significance

p.f. [2] Z_{α} we reject H_0
 and conclude that there is a
 significant difference.
 If $|Z| < Z_{\alpha}$ we accept H_0
 and conclude that there is no
 significant difference.

Test of significance for large samples:

\Rightarrow (n > 30) then large sample

\Rightarrow Classified as:

- 1) Testing of significance for ~~single~~ single proportion:
- 2) Testing of significance for different proportion.
- 3) Testing of significance for single mean.
- 4) Testing of significance for different means.

1) For single proportion:

Let X be the no. of successes in

n independent trials with

constant probability 'p' of success

of each trial

$$E(X) = np$$

$$V(X) = npq$$

$Q = 1 - p$
 probability of failure

$p = X/n$ called observed proportion of success.

$$E(p) = E(X/n) = \frac{1}{n} E(X) = \frac{np}{n} = p$$

$$V(p) = V(X/n) = \frac{1}{n^2} V(X) = \frac{npq}{n^2} = \frac{pq}{n}$$

$$S.E(p) = \sqrt{\frac{pq}{n}}$$

$$Z = \frac{p - E(p)}{S.E(p)} = \frac{p - p}{\sqrt{\frac{pq}{n}}} = 0$$

$$Z = \frac{p - p}{\sqrt{\frac{pq}{n}}} = 0$$

Z is called test static

Note:

1) The probable limit for the observed proportion of success are

$$p \pm 2 \sqrt{\frac{pq}{n}}$$

q is 1-p to the level of significance

2) If small 'p' is not known, the limits for the proportion in the population are

$$p \pm z_{\alpha} \sqrt{\frac{pq}{n}}$$

$$q = 1 - p$$

$$p \pm z_{\alpha} \sqrt{\frac{pq}{n}}$$

$$q = 1 - p$$

3) If α is not given so we can take safely 3% limits.

4) Confidence limits for the observed proportion

$$p \pm 3 \sqrt{\frac{pq}{n}}$$

5) Confidence limits for the population proportion

$$p \pm \sqrt{\frac{pq}{n}}$$

eg) A coin was tossed 400 times and the head turned up 216 times

Test the hypothesis that the coin is unbiased.

Ans) H_0 = The coin is unbiased

$$p = \frac{1}{2} = 0.5$$

H_1 = The country is not unbiased (biased)

$P \neq 1/2 = 0.5$

$n = 400$

$X = \text{no. of supporters} = 216$

$p(X) = \frac{216}{400} = 0.54$

P = population proportion = 0.5

$Q = 1 - P = 0.5$

Test statistic

$|Z| = \frac{P - P_0}{\sqrt{P_0 Q / n}}$

$= \frac{0.54 - 0.5}{\sqrt{\frac{0.5 \times 0.5}{400}}}$

$Z = 1.6$

Since $Z < Z_{\alpha}$ we use the tailed test

Conclusion is:

Since $|Z| = 1.6 < 1.96$

$|Z| < Z_{\alpha}$ is the significant value at 5% level of significance

27 A certain suspect die has thrown 9000 times and 5 or 6 was obtained 3240 times. On the assumption of certainty throwing of the plate indicate an unbiased die.

$$H_0 = 2/6 = 1/3 \quad F.P.$$

$$g = 2/3 \quad N = 9000$$

$$F = \frac{3240}{9000} = 0.36$$

$$P = 0.36$$

$$H_1 \neq 1/3 \quad P$$

$$\frac{2}{6} = 1/3$$

Two tailed test

$$|Z| \leq \frac{0.36 - 1/3}{\sqrt{\frac{1/3 \times 2/3}{9000}}} < |Z|$$

$$\sqrt{\frac{1/3 \times 2/3}{9000}}$$

$$= 0.005$$

$$= 6009 \quad 4.2 \times 10^{-4}$$

Hypothesis is rejected as

$$|Z| > Z_\alpha$$

37 A manufacturer claims that only 4% of product supplied by him are defective. A random sample of 600 products containing

36 defectives. Test the claim of manufacturer.

$$0.04 = 4\%$$

$$P = 0.06$$

$$Q = 0.94$$

$$p = 0.04$$

$$H_0 = 0.04$$

$$Z = \frac{p - P}{\sqrt{PQ/n}}$$

$$= \frac{0.04 - 0.06}{\sqrt{\frac{0.06 \times 0.94}{600}}}$$

$$= \frac{-0.02}{\sqrt{0.00094}}$$

$$= -2.0261$$

$$|Z| > |Z_{\alpha/2}| \text{ so two tailed}$$

test -

so, manufacturer's claim is

rejected.

Test of difference b/w proportions:

Consider two samples X_1 and X_2

of size n_1 and n_2 respectively

taken from two different populations

To test the significance of the difference b/w sample proportions

P_1 and P_2 . The test statistic

under the null hypothesis H_0 ,
that there is no significant
difference b/w 2 sample proportions

$$|Z| = \frac{P_1 - P_2}{\sqrt{PQ \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

$$P = \frac{n_1 P_1 + n_2 P_2}{n_1 + n_2}$$

$$Q = 1 - P$$

eg) A machine produced 16 defective
articles in a batch of 500

After overally it produces 3
defectives in a batch of 100.
Has the machine improved?

$$n_1 = 500, n_2 = 100$$

$$P_1 = \frac{16}{500} = 0.032, P_2 = \frac{3}{100} = 0.03$$

Null hypothesis H_0 =

Alternative hypothesis H_1 = $P_1 > P_2$

(right tail)

$$P = \frac{n_1 P_1 + n_2 P_2}{n_1 + n_2} = 0.031$$

$$Q = 1 - 0.031 = 0.969$$

$$Z = \frac{0.082 - 0.031}{\sqrt{0.031 \times 0.968 \left(\frac{1}{500} + \frac{1}{500} \right)}}$$

$$= \frac{0.105}{0.01714}$$

$$\text{Conclusion: } \frac{0.105}{0.01714} = 6.12$$

$$|Z| < 1.645 \quad (95\%)$$

level of significance less than

$$5\%$$

accepted - is not improved machine

substantially than accepted
- ~~not~~ greater than rejected

2% good articles from a factory are examined and found to be

2% defective, can it reasonably be concluded that the factory is

of the first factory are 2% defective to the 2% of

$$n_1 = 500, n_2 = 800$$

$$EOP = 0.021 - 0.015$$

$$\sqrt{\frac{0.021 \times 0.979}{500} + \frac{0.015 \times 0.985}{800}}$$

$$|Z| =$$

$$P = \frac{P_1 n_1 + P_2 n_2}{n_1 + n_2}$$

$$Z = \frac{10 + 12}{1300}$$

$$= 0.016$$

(complete later)

of significance for
test for single mean

i.e. to test whether difference
between sample mean and population
mean is significant or not.

$X_1, X_2, X_3, \dots, X_N$ are

sample size n from large
population

X_1, X_2, \dots, X_N are
with mean μ and variance σ^2

The standard error of mean
of random sample of
size n from population
with variance σ^2 is

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

$$|z| = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

σ : standard deviation.

if σ is not known,

$$|z| = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

Note:

at the level of significance α and n is the critical value then

$$|z| < |z_{\alpha}| \Rightarrow \left| \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \right| < z_{\alpha}$$

limit of the population mean

$$\bar{x} - z_{\alpha} \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + z_{\alpha} \frac{\sigma}{\sqrt{n}}$$

At 5% level of significance, 95% confidence limit are:

$$\bar{x} - 1.96 \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + 1.96 \frac{\sigma}{\sqrt{n}}$$

At 1% level of significance, 99% confidence limit are given by:

$$\bar{x} - 2.58 \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + 2.58 \frac{\sigma}{\sqrt{n}}$$

These limits are called confidence limits.

A Normal population has a mean of 6.8 and standard deviation of 1.5. A sample of 400 members gave a mean of 6.75, is the difference significant?

Given:

$$\mu = 6.8, \sigma = 1.5, n = 400$$

$$\bar{x} = 6.75$$

$$Z = \frac{6.75 - 6.80}{1.5 / \sqrt{400}}$$

$$= \frac{-0.05}{0.075} = -0.67$$

H_0 there is no significant difference b/w the two.

H_1 there is significant difference b/w \bar{x} and μ

$$|Z| = 0.67$$

Conclusion:

$$|Z| < Z_{\alpha} = 1.96 \quad \text{5\% level of significance}$$

Selected H_0 is accepted.

Q7 The mean weight obtained from a random sample of size 100 is 64 gms. The standard deviation of the population is 3 gms. Test the weight of the population is 64 gms at 5% level of significance.

Also setup H_0 & H_1 and write the results of the mean weight of the population.

$$n = 100, \mu = 64, \sigma = 3, \alpha = 5\%$$

H_0 : There is no significant difference b/w sample & population mean.

$$H_1: \mu \neq 64$$

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{64 - 64}{3 / \sqrt{100}} = 0$$

$$|Z| > 1.96$$

Conclusion: $\mu = 64$ is not significant.

Conclusion: H_0 is rejected.
 i.e. Sample is not drawn from the population with $\mu = 67$.
 To find confidence limit.

$$\bar{x} \pm z \cdot \frac{\sigma}{\sqrt{n}}$$

$$= 64.3229 \pm 2.226 \cdot \frac{6}{\sqrt{10}}$$

} any one of true mean

Test of significance for difference of means of 2 large samples.
 Let \bar{x}_1 be the mean of sample size n_1 , from a population with mean μ_1 , and variance σ_1^2 .

Similarly \bar{x}_2 be the mean of sample size n_2 , from a population mean μ_2 and variance σ_2^2 .

$$|Z| = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

Under null hypothesis samples are drawn from the sample.

population, where $\sigma_1 = \sigma_2$
and $\mu_1 = \mu_2$

$$Z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma^2}{n_1} + \frac{\sigma^2}{n_2}}}$$

$$Z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma^2 + \sigma^2}{n_1 + n_2}}}$$

Note: σ_1 and σ_2 are not known
or $\sigma_1 \neq \sigma_2$

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

or σ is not known and

$$\sigma_1 = \sigma_2$$

$$\sigma^2 = \frac{1}{n_1 + n_2} (s_1^2 + s_2^2)$$

$$\sigma = \sqrt{\frac{1}{n_1 + n_2} (s_1^2 + s_2^2)}$$

Ques) The average income of a person was of ₹ 210 with a standard deviation of

₹ 10, is sample of 100 people of a city. For another sample of 150 persons, the average was ₹ 220 with standard deviation of ₹ 12. The standard deviations of income of people of the city was significant whether there is any difference between average incomes of the localities.

$$n_1 = 100$$

$$\sigma_1 = 10$$

$$n_2 = 150$$

$$\sigma_2 = 12$$

$$H_0 =$$

$$\mu_1 = \mu_2$$

$$\mu_1 \neq \mu_2$$

$$H_1 = \mu_1 \neq \mu_2$$

$$Z = \frac{210 - 220}{\sqrt{\frac{100}{100} + \frac{144}{150}}}$$

$$= -1.4$$

$$70.14$$

Conclusion: H_0 at 5% level

is Rejected at significant

Ques 2) For sample 1 $N_1 = 11,210$

$$\sigma_x^2 = 49,000$$

$$(\sigma_x - \bar{x})^2 = 7,84,000$$

For sample 2

$$N_2 = 1500$$

$$\sigma_x^2 = 70,500$$

$$(\sigma_x - \bar{x})^2 = 24,00,000$$

Discuss the significance of the difference of sample mean.

Null hypothesis H_0 of no significant difference

H_1 = $\bar{x}_1 \neq \bar{x}_2$

To calculate the

$$\bar{x}_1 = \frac{7,84,000}{11,210}$$

$$\bar{x}_2 = \frac{24,00,000}{1500}$$

$$\bar{x}_1 = 70,000$$

$$\bar{x}_2 = 16,000$$

$$\bar{x}_1 = \frac{629}{n_1} = 49$$

$$\bar{x}_2 = 47$$

$$s_1 = 10.49$$

$$s_2 = \sqrt{\frac{384^2}{1000} + \frac{1600^2}{1500}}$$

10% Accepted

Test of significance for difference of standard deviations

$$Z = \frac{s_1 - s_2}{\sqrt{\frac{s_1^2}{2n_1} + \frac{s_2^2}{2n_2}}}$$

σ_1 and σ_2 are Population

Standard deviation. When population standard deviation are not known,

$$Z = \frac{s_1 - s_2}{\sqrt{\frac{s_1^2}{2n_1} + \frac{s_2^2}{2n_2}}}$$

$$\sqrt{\frac{384^2}{1000} + \frac{1600^2}{1500}}$$

Q7 Random sample and data from 2 countries gave the following data relating to the height of the adult males.

	Country A	Country B
Mean height	67.42	67.25
S.D	2.58	2.50
Number of sample	1000	200

1) Is the difference b/w the means as significant?

2) Is the difference b/w the means significant?

Soln) $n_1 = 1000$, $n_2 = 200$

$\bar{x}_1 = 67.42$, $\bar{x}_2 = 67.25$

$s_1 = 2.58$, $s_2 = 2.50$

Since sample size are large we can take

$$\bar{G}_1 = 20.58$$

$$\bar{G}_2 = 20.50$$

$$Z = \frac{\bar{G}_1 - \bar{G}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{20.58 - 20.50}{\sqrt{\frac{0.00011}{100} + \frac{0.00011}{100}}} = 1.56$$

$$|Z| < 1.96$$

Accepted null hypothesis

Accepted for 5% level of significance

Level of significance

Not significant

not

to 0.01

28	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	-----

Chi square test for goodness of fit

Problem

$$\sum_{i=1}^n \left(\frac{(O_i - E_i)^2}{E_i} \right)$$

The following table gives the no. of accidents that took place in an industry during various days of a week.

Test if accidents are uniformly distributed over the week.

Day	No. of accidents
No. of accidents	

Day	M	T	W	Th	fr	Sa
No. of accidents	14	18	12	11	15	14
					84	

H_0 = Accidents are uniformly distributed over the week.

Expected frequency of the accidents on each days of the week = 84
total days = 6

expected = 14 accidents

observed frequency

	14	18	12	11	15	14
E_i	14	14	14	14	14	14

$$(0)^2 + 16 + 4 + 9 + 1 + 0$$

$$\frac{30}{14} = \frac{15}{7} = 2.14$$

χ^2
D.O.F = 5
 $n = 6$

$n = 6$ D.O.F = 5

$$\chi^2 \text{ at } 5\% = 11.49$$

Test from D.O.F table
if below 5% significance
then accepted else vice versa.

Records taken of the no. of male and female births in 800 families having 4 children.

No. of male birth	No. of female birth	
0	4	32
1	3	178
2	2	290
3	1	236
4	0	94
		0

Test whether the data are consistent with the hypothesis that the binomial law holds and the chance of male birth is equal to that of female birth.

$$n = 800, p = q = \frac{1}{2}$$

$$x = 0, 1, 2, 3, 4$$

$$P(X=x) = \binom{n}{x} p^x q^{n-x}$$

$$H_0: p = q = \frac{1}{2}$$

$$H_1: p \neq q$$

(N=) total frequency

N₀ = No. of families with 0 male children.

$$P(X) = \frac{N_0}{N}$$

$$P(X=x) = n C_x p^x q^{n-x}$$

$$x = 0, 1, 2, 3, 4$$

$$N = 830$$

$$H(x) = N \cdot P(X=x)$$

$$N(0) = 830 \times {}^n C_0 \left(\frac{1}{2}\right)^0 \left(\frac{1}{2}\right)^4 = 51.8$$

$$N(1) = 830 \times {}^n C_1 \left(\frac{1}{2}\right)^1 \left(\frac{1}{2}\right)^3 = 207.5$$

$$N(2) = 830 \times {}^n C_2 \left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right)^2 = 51.8$$

$$N(3) = 830 \times {}^n C_3 \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^1 = 207.5$$

$$N(4) = 830 \times {}^n C_4 \left(\frac{1}{2}\right)^4 \left(\frac{1}{2}\right)^0 = 51.8$$

$$51.8, 207.5, 51.8, 207.5, 51.8$$

$$\frac{(32-51.8)^2}{51.8} + \frac{(178-207.5)^2}{207.5} + \frac{(99-51.8)^2}{51.8}$$

$$+ \frac{(236-207.5)^2}{207.5} + \frac{(99-51.8)^2}{51.8}$$

$$= 7.5 + 9.19 + 1.48 + 3.91$$

$$= 31.39$$

$$= 51.459$$

$$D.O.F =$$

$$9.49$$

Rejected

$$\frac{(32-51.8)^2}{51.8} + \frac{(178-207.5)^2}{207.5} + \frac{(99-51.8)^2}{51.8}$$

