

eg:-

Annual truck repair expenses * Repair expense during last yr in hundred of \$

Truck no	Age of truck in yrs	Repair expense during last yr in hundred of \$
201	5	7
102	3	7
103	3	6
104	1	4

sol:-

$$b = \frac{\sum xy - n \cdot \bar{x} \cdot \bar{y}}{\sum x^2 - n \bar{x}^2} = \frac{(35 + 21 + 18 + 4) - 4 \times 3 \times 6}{(25 + 9 + 9 + 1) - 4 \times 9}$$

$$= \frac{68 - 72}{44 - 36} = \frac{-4}{8} = -0.5$$

$$a = \bar{y} - b \cdot \bar{x} = 6 - 0.5 \times 3 = 3.75$$

Reg. eqn $\Rightarrow y = 3.75 + 0.25x$

If the age of truck is 6 then.

estimated value (repair expense) is $\hat{y} = 3.75 + 6 \times 0.25$

Standard error

$$S_e = \sqrt{\frac{\sum y - \sum \hat{y}}{n-2}}$$

y = act. value, \hat{y} = estimated value

When we have regression eqn, we will have 'n' pairs of datasets, the regression eqn $y = a + bx$ and a and b are reg. coefficients and are estimated values. So, degrees of freedom $n-2$. because the 'n' no. of

data points are linked by 2 normal equations. i.e. 2 equations.

So, degree of freedom is $n-2$.
Standard Error is measure of reliability of regression eqn.

If SE is higher i.e. dist b/w \hat{y} and \hat{y} (reg line) then points are more dispersed / they are far away.

If SE is 0 i.e. all the points are on a line. Then we can say that regression line is perfect fit for data points.

Standard Error for some problem, both can be used.

$$S_e = \sqrt{\frac{\sum(Y-\hat{y})^2}{n-2}} \quad S_e = \sqrt{\frac{\sum Y^2 - 2\sum XY + \sum X^2}{n-2}}$$

$$S_b = \frac{S_e}{\sqrt{\sum X^2 - n\bar{X}^2}}$$

$$\hat{y} = 3.75 + 0.75x \Rightarrow \text{if } x = 5 \Rightarrow \hat{y} = 3.75 + 3.75 = 7.5$$

$$n = 3 \Rightarrow \hat{y} = 3.75 + 2.25 = 6$$

$$S_e = \sqrt{\frac{6.25 + 0.25}{4-2}}$$

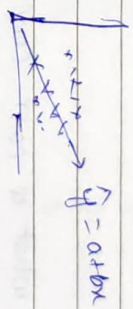
$$n = 1 \Rightarrow \hat{y} = 3.75 + 0.75 = 4.5$$

$$S_e = \sqrt{\frac{(1-3.75)^2 + (1-6)^2 + (6-4.5)^2 + (4-4.5)^2}{4-2}}$$

$$= \sqrt{\frac{0.25 + 1 + 0 + 0.25}{2}} = \sqrt{\frac{1.5}{2}} = \sqrt{0.75} = 0.866$$

There is standard error of equation.

If $S_e = 0$ then we would have said line is passing through points. But here value is 0.866 i.e. some points are away from line.



In a normal distribution, there are 68% points in ± 1 SD, 95% in ± 2 SD.

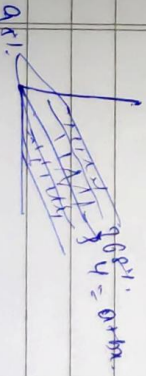
If the size of sampling dist > 30 then 68% of sample points then 68% in sample mean ± 1 SD.

95% points lie in sample mean ± 2 SD.

$$\pm 2 = \frac{\bar{x} - \mu}{S.E.} \quad \mu \pm 3 S.E. \rightarrow 99\%$$

$$\hat{N} = \bar{x} \pm Z \cdot S.E. \quad 95\% - 99\%$$

99% confidence level is wider than 95%.



when we have large data set then we can take t_{90} follows Z -distribution.

Conf. interval = $(\hat{q} \pm t_{90} SE)$

For smaller sample size t_k follows t distribution.

Conf. interval = $(\hat{q} \pm t_{k} SE)$

if age of truck = u yrs. what is corr. conf. int at 90%

$\hat{q} = 8.75 + u \times 0.75 = 8.75$ as we have only u observed we used t -distribution.

$= 8.75 \pm t_{90} \times 0.866 \Rightarrow t_{90}$ at $DOF = 2 \Rightarrow$

$= [6.75 - 2.920 \times 0.866, 8.75 + 2.920 \times 0.866]$
 $[4.22, 9.22]$ based on data that we have

For a truck of age u yrs, with 90% confidence we can say that repair expense would be b/w (4.22, 9.22)

Here, we also have the standard error for b-coefficient in $y = a + bx$.

$\Rightarrow S_b = \frac{SE}{\sqrt{\sum X^2 - n\bar{X}^2}}$

eg:-

Yr	Bad Exp	Profit
2013	5	31
2014	11	40
2015	4	30
2016	5	34
2017	3	25
2018	2	20

pdf:-

y can be expressed as $y = a + bx$

$SE = \sqrt{\frac{y - \hat{y}}{n-2}}$
 $b = \frac{\sum xy - n\bar{x}\bar{y}}{\sum x^2 - n\bar{x}^2}$

$\hat{y} = a + b \cdot x = 30.5$
 $b = \frac{(1035 + 440 + 120 + 120 + 25 + 40) - 6 \times 30 \times 5}{(2.5 + 121 + 16 + 25 + 9 + 4) - 6 \times 25}$

$b = \frac{1000 - 900}{200 - 150} = \frac{100}{50} = 2$

$\hat{y} = 20 + 2x$
 $\hat{y} = 20 + 2 \times 11 = 42$

$SE = \sqrt{\frac{(1) + 4 + 4 + 16 + 1 + 16}{4}}$
 $x = 11 \Rightarrow 20 + 2 \times 11 = 42$

$= \sqrt{42/4} = \sqrt{10.5}$
 $= 3.24$

$$S_b = \frac{S_e}{\sqrt{\sum x_i^2 - n\bar{x}^2}}$$

$$= \frac{3.24}{\sqrt{(95+121+18+19+10) - 6 \cdot (9)^2}}$$

$$= \frac{3.24}{\sqrt{200 - 150}} = \frac{3.24}{\sqrt{50}} = \frac{3.24}{7.07} = 0.458$$

S_b is dependent on S_e .

Conf. int for b coeff at 90% conf. level.

$$\text{conf. int} = (\hat{b} \pm t_{\alpha/2} \cdot S_b)$$

$$= (2 \pm 2.132 \cdot 0.458)$$

$$= (2 \pm 0.98) = (1.02, 2.98)$$

at 90% conf. level.

if b value i.e., slope.

For every \uparrow in R&D exp profit will be increasing (102, 298)

if we have less conf. level the range will be less.

For simple regression, it depends on 2 normal eqns.

If we consider S_e as a whole then it will effect whole \uparrow

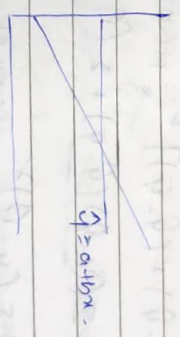
based on, Reg. coeff 'b' have can have conf. int. every \uparrow in x value what will be the \uparrow in y .

01/12/21

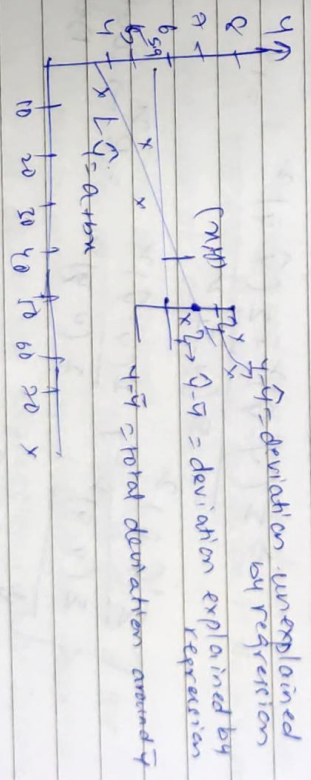
For a certain value of independent variable the confidence interval of dependent variable can be determined using sample error of equation.

Explaining variation: R^2

A. Breaking down the distances



Total Variation:



Total Deviation = Explained Deviation + Unexplained Deviation

explained \rightarrow when the point is on the line i.e., it satisfies the regression equation. So, it been explained by regression.
 unexplained \rightarrow If the line is not passed through point i.e., not explained by regression i.e., arbitrary point.

Sums of Squares

We can sum this equation across all the y's and square both sides to get:

$$\begin{aligned} \sum (y - \hat{y})^2 &= \sum [(y - \hat{y}) + (\hat{y} - \bar{y})]^2 && \text{P.C.T. a theorem that} \\ &= \sum (y - \hat{y})^2 + 2 \sum (y - \hat{y})(\hat{y} - \bar{y}) + \sum (\hat{y} - \bar{y})^2 && \text{prod. 2} \\ &= \sum (\hat{y} - \bar{y})^2 + \sum (y - \hat{y})^2 && \text{Square of } \end{aligned}$$

$\sum (y - \hat{y})^2$ Square of exp. deviation
 $\sum (\hat{y} - \bar{y})^2$ Square of unexplained deviation

$$\begin{aligned} \sum (y - \bar{y})^2 &= \sum (y - \hat{y})^2 + \sum (\hat{y} - \bar{y})^2 \\ \Rightarrow \frac{\sum (y - \hat{y})^2 + \sum (\hat{y} - \bar{y})^2}{\sum (y - \bar{y})^2} &= 1 \end{aligned}$$

$$\frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2} = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2} = R^2 \rightarrow \text{Coefficient of Determination}$$

Rel. Deviation

How many % of points of plane can be explained by regression eq. in scatter plot that is those many points lie on line.

if $R^2 = 0.81$
 $= \sqrt{0.81} = \pm 0.9$

i.e., it can be +ve or -vely correlated.

It depends on the slope of the line.

If the slope is +ve then the coefficient is +ve,
 " " " " " " " " -ve
 " " " " " " " " -ve

Variation of y values around the regression line:

$$\sum (y - \hat{y})^2$$

Variation of y values around their own values

$$\sum (y - \bar{y})^2$$

Coefficient of determination r^2 .

Total sum of squares (CST)

the term on the left hand side of this eqn is

the sum of the squared distances from all points to \bar{y} . we call this the total variation in the y's, or

the total sum of squares (CST).

Regression sum of squares:

the first term on the right hand side is the

sum of the squared distances from the regression

line to \bar{y} . we call it the Regression sum of squares,

or CSR.

Error sum of squares

Finally, the last term is the sum of squared distances

from the point to the regression line. Remember, this is the quantity that least squares minimizes.

$$TSS = SSE + SSE$$

P-Value Approach

The area under the density function right to the value of Z obtained from the sample,

when $P < \alpha$

Alternate/research hypothesis gets accepted at α .

Significance level.

When $P < \text{level of sig}$ then alt. hypothesis gets accepted.

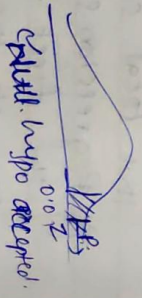
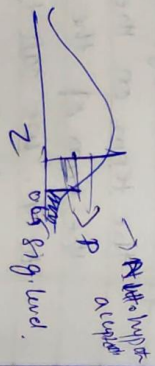
" $P > \alpha$ " " " " " " " " " " " "

Instead of critical value we are checking P value.

$$P_{1-\alpha}$$

P is right side of Z

P is compared to significance level.



Multiple regression

$$Y = a + b_1 X_1 + b_2 X_2$$

$$SE = \frac{\sum(Y - \hat{Y})^2}{n-2}$$

Y \rightarrow dep. var
 $X_1, X_2 \rightarrow$ indep variables

Y = a + bX
 $SE = \frac{\sum(Y - \hat{Y})^2}{n-2}$ single regression

$$SE = \sqrt{\frac{\sum(Y - \hat{Y})^2}{(n-k-1)}} \quad k = \text{no. of independent variable}$$

In this case $k=2 \Rightarrow n-2-1$

In prev. case indep var = 1 $\Rightarrow k=1 \Rightarrow n-1-1 = n-2$

$$TSS = SSD + SSE$$

Similar to F-test in regression also we have 2 variances i.e. SSE and SSD.

We can also do a F-test to find significance of regression as a whole.

$$F = \frac{SSD/k}{SSE/(n-k-1)}$$

If $P > F$ then we can say reg. eqn is significant.

Pop. reg. eqn $Y = A + B_1 X_1 + B_2 X_2$

$$H_0: B_1 = B_2 = 0$$

$$H_1: B_1 \neq B_2 \neq 0$$

If null hypo is true X_1 and X_2 cannot explain regression as whole.

When alt. hypothesis is true then we can say reg. can be explained in terms of independent variables.

$$H_0: B_1 = 0 \quad H_1: B_1 \neq 0$$

This is regarding one coefficient B.
 - then we use t-test.

$$t = \frac{B_1 - 0}{SE_B}$$

with the help of t-test we can tell for individual coefficients.

ANOVA in Excel

MS \rightarrow Mean Squares

SS \rightarrow sum of squares.

in reg,

$$\text{Adjusted } R^2 = 1 - (1 - R^2) \left(\frac{N-1}{N-k-1} \right)$$

^{No. of}
 $N = \text{Data points}$ $k = \text{No. of independent variables}$

when k is small when compared to 1 then $\frac{N-1}{N-k-1} \approx 1$.

$$\Rightarrow \text{Adjusted } R^2 = R^2$$