

Linear Regression and Correlation

The simple linear regression consider a single regressor variable or predictor variable X and a dependent or response variable Y .

$$Y = \beta_0 + \beta_1 X + \epsilon \quad - (1)$$

Where, the intercept β_0 and the slope β_1 are unknown regression coefficients. ϵ is a random error.

Least Square Method-

The method of least square uses to estimate the regression coefficient.

Using equation (1), we may express the n observations in the sample as

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad i=1, 2, \dots, n$$

and the sum of the squares of the individual deviations of the observations from the true regression line is

$$L = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

The least squares estimators of β_0 and β_1 , say $\hat{\beta}_0$ and $\hat{\beta}_1$ must satisfy

$$\left. \frac{\partial L}{\partial \beta_0} \right|_{\hat{\beta}_0, \hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

$$\left. \frac{\partial L}{\partial \beta_1} \right|_{\hat{\beta}_0, \hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0$$

Simplifying the above two equation gives

$$\left. \begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n y_i x_i \end{aligned} \right\} (2)$$

Equations (2) are called the least square normal equations. The solution to the normal eqn results in the least squares estimators $\hat{\beta}_0$ and $\hat{\beta}_1$.

Least Squares Estimates-

The least squares estimates of the intercept and slope in the simple linear regression model are

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$
$$\hat{\beta}_1 = \frac{\sum_{i=1}^n y_i x_i - \frac{(\sum_{i=1}^n y_i)(\sum_{i=1}^n x_i)}{n}}{\sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}}$$

$$\text{where } \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \text{ and } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

The fitted estimated regression line -

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

Q find the line of regression on the basis of the following data

X:	1	2	3	4	5	6	7
Y:	2	4	7	6	5	6	5

Soln

X	Y	X^2	Y^2	XY
1	2	1	4	2
2	4	4	16	8
3	7	9	49	21
4	6	16	36	24
5	5	25	25	25
6	6	36	36	36
7	5	49	25	35
$\Sigma X = 28$	$\Sigma Y = 35$	$\Sigma X^2 = 140$	$\Sigma Y^2 = 191$	$\Sigma XY = 151$

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^7 Y_i = \frac{1}{7} \times 35 = 5$$

$$\bar{X} = \frac{1}{n} \sum_{i=1}^7 X_i = \frac{1}{7} \times 28 = 4$$

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum Y_i X_i - \frac{(\sum Y_i)(\sum X_i)}{n}}{\sum X_i^2 - \frac{(\sum X_i)^2}{n}} \\ &= \frac{151 - \frac{(35)(28)}{7}}{140 - \frac{(28)^2}{7}} \\ &= \frac{151 - 140}{140 - 112} \\ &= \frac{11}{28} \end{aligned}$$

$$\begin{aligned} \hat{\beta}_0 &= \bar{Y} - \hat{\beta}_1 \bar{X} \\ &= 5 - \frac{11}{28} \times 4 \\ &= 5 - \frac{44}{28} \\ &= \frac{140 - 44}{28} \\ &= \frac{96}{28} = \frac{24}{7} \end{aligned}$$

The regression line is

$$\hat{Y} = \frac{24}{7} + \frac{11}{28} X \quad \text{or} \quad \hat{Y} = 3.42 + 0.39X$$

Q Fit a least square quadratic fit to the following data set.

X :	0	1	2	3	4
Y :	1.0	1.5	1.5	2.5	3.5

Soln We have to fit the quadratic fit $y = a + bx + cx^2$

The three normal eqn for determining the value of numerical constants a, b and c are :

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

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

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

X	Y	x^2	x^3	x^4	xy	x^2y
0	1.0	0	0	0	0	0
1	1.5	1	1	1	1.5	1.5
2	1.5	4	8	16	3	6.0
3	2.5	9	27	81	7.5	22.5
4	3.5	16	64	256	14	56.0
10	10	30	100	354	26	86

Substituting the values obtained from the in the normal eqn

$$10 = 5a + 10b + 30c \quad - (1)$$

$$26 = 10a + 30b + 100c \quad - (2)$$

$$86 = 30a + 100b + 354c \quad - (3)$$

Multiplying eqn (1) by 2 and then subtract from eqn (2) we get

$$6 = 10b + 40c \quad - (4)$$

Next multiplying eqn (2) by 3 and then subtracting from eqn (3)

$$8 = 10b + 54c \quad - (5)$$

Now subtracting eqnⁿ (4) from eqnⁿ (5) we get

$$2 = 14C$$

$$\Rightarrow C = \frac{2}{14} = \frac{1}{7} = 0.143$$

Substituting $C = \frac{1}{7}$ in eqnⁿ (4) we get

$$6 = 10b + 40 \times \frac{1}{7}$$

$$b = 0.0286$$

Substituting the values of b and C in eqnⁿ (1) we get

$$10 = 5a + \frac{10}{35} + \frac{30}{7}$$

$$\Rightarrow 5a = 10 - \frac{2}{7} - \frac{30}{7}$$

$$\Rightarrow a = \frac{38}{35} = 1.0857$$

Thus, $a = 1.0857$, $b = 0.0286$ and $C = 0.143$

Hence, the non-linear fitted line is

$$\hat{y} = 1.0857 + 0.0286x + 0.143x^2$$

Regression On Transformed Variable -

We occasionally find that the straight-line regression model $Y = \beta_0 + \beta_1 X + \epsilon$ is inappropriate because the true regression function is nonlinear. In some of these situations, a nonlinear function can be expressed as a straight line by using a suitable transformation. Such nonlinear models are called intrinsically linear.

As an example of a nonlinear model that is intrinsically linear, consider the exponential function

$$Y = \beta_0 e^{\beta_1 X} + \epsilon$$

The function is intrinsically linear, since it can be transformed to a straight line by a logarithmic transformation.

$$\log Y = \log \beta_0 + \beta_1 X + \log \epsilon$$

Q Fit an exponential curve of the form $Y = Ae^{BX}$ for the following data:

X :	1	2	3	4
Y :	7	11	17	27

Note -

1. Fitting of a power Curve $Y = ax^b$ to a set of n points,
 (i)

Taking logarithm of each side, we get

$$\log Y = \log a + b \log X$$

$$\Rightarrow U = A + bV$$

Where, $U = \log Y$, $A = \log a$ and $V = \log X$

This is a linear eqn in V and U

Normal equations for estimating A and B are;

$$\sum U = nA + b \sum V$$

$$\text{and } \sum UV = A \sum V + b \sum V^2$$

These equations can be solved for A and b and consequently, we get

$a = \text{antilog}(A)$.

With the values of ' a ' and ' b ' so obtained (i) is the curve of best fit to the set of n points.

2. Fitting of Exponential Curves (i) $Y = ab^x$, (ii) $a e^{bx}$ to a set of n points.

(i) $Y = ab^x$

Taking logarithm of each side, we get

$$\log Y = \log a + x \log b$$

$$\Rightarrow U = A + Bx$$

Where $U = \log Y$, $A = \log a$ and $B = \log b$

This is a linear eqn in x and U

The normal equations for estimating A and B are:

$$\sum U = nA + B \sum X$$

$$\text{and } \sum XU = A \sum X + B \sum X^2$$

Solving these equations for A and B, we finally get

$$a = \text{Antilog}(A) \text{ and } b = \text{Antilog}(B)$$

With these values of 'a' and 'b' (i) is the curve of the best fit to the given set of n points.

(ii) $Y = ae^{bx}$

$$\begin{aligned} \log Y &= \log a + bx \log e \\ &= \log a + (b \log e) x \end{aligned}$$

$$\Rightarrow U = A + BX$$

$$\text{Where, } U = \log Y, A = \log a \text{ and } B = b \log e$$

This is a linear equation in X and U, and the normal equations are:

$$\sum U = nA + B \sum X$$

$$\text{and } \sum XU = A \sum X + B \sum X^2$$

From these we find A and B consequently $a = \text{Antilog } A$ and

$$b = \frac{B}{\log e}$$

Soln Taking logarithm of each side we get

$$\log Y = \log A + B \log e$$

$$\Rightarrow U = a + bx$$

Where, $U = \log Y$, $a = \log A$ and $b = B \log e$

This is a linear equation in x and U and the normal equations are:

$$\begin{aligned} \sum U &= na + b \sum x \\ \sum xU &= a \sum x + b \sum x^2 \end{aligned} \quad \text{--- (1)}$$

X	Y	U = log Y	XU	X ²
1	7	0.845	0.845	1
2	11	1.041	2.082	4
3	17	1.230	3.690	9
4	27	1.431	5.724	16
<u>10</u>		<u>4.547</u>	<u>13.341</u>	<u>30</u>

Substitute these values in eqn (1) we get

$$4.547 = 4a + 10b \quad \text{--- (3)}$$

$$13.341 = 10a + 30b \quad \text{--- (4)}$$

After solving eqn (3) and (4) we get

$$a = 0.15 \quad \text{and} \quad b = 0.3947$$

$$A = \text{antilog}(a) = \text{antilog}(0.15) = 1.413$$

$$B = \frac{b}{\log e} = \frac{0.3947}{0.4342} = 0.9090$$

The fitting curve is $Y = 1.413 e^{(0.9090)x}$

Correlation-

If the change in one variable affects a change in the other variable, the variables are said to be correlated.

If the two variables deviate in the same direction, i.e., if the increase (or decrease) in one results in a corresponding increase (or decrease) in the other, correlation is said to be direct or positive.

If the two variables deviate in the opposite direction i.e., if the increase (or decrease) in one results in a corresponding decrease (or increase) in the other, correlation is said to be diverse or negative.

For example - The income and expenditure is positively correlated.
The price and demand is negatively correlated.

Karl Pearson's coefficient of correlation-

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

$$\text{Where, } \text{Cor}(x, y) = \frac{1}{n} \sum x_i y_i - \bar{x} \bar{y}$$

$$\sigma_x^2 = \frac{1}{n} \sum x_i^2 - \bar{x}^2$$

$$\sigma_y^2 = \frac{1}{n} \sum y_i^2 - \bar{y}^2$$

Note:- Correlation coefficient always lies between -1 and +1.
If $r = +1$, the correlation is perfect and positive.
If $r = -1$, the correlation is perfect and negative.

Q Calculate the correlation coefficient for the following heights (in inches) of fathers (X) and their sons (Y):

X:	65	66	67	67	68	69	70	72
Y:	67	68	65	68	72	72	69	71

Soln

X	Y	X ²	Y ²	X Y
65	67	4225	4489	4355
66	68	4356	4624	4488
67	65	4489	4225	4355
67	68	4489	4624	4556
68	72	4624	5184	4896
69	72	4761	5184	4968
70	69	4900	4761	4830
72	71	5184	5041	5112
<u>544</u>	<u>552</u>	<u>37028</u>	<u>38132</u>	<u>37560</u>

$$\bar{X} = \frac{1}{n} \sum X = \frac{544}{8} = 68, \quad \bar{Y} = \frac{1}{n} \sum Y = \frac{552}{8} = 69$$

$$\begin{aligned} r(X, Y) &= \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\frac{1}{n} \sum XY - \bar{X} \bar{Y}}{\sqrt{\left(\frac{1}{n} \sum X^2 - \bar{X}^2\right) \left(\frac{1}{n} \sum Y^2 - \bar{Y}^2\right)}} \\ &= \frac{\frac{1}{8} \times 37560 - 68 \times 69}{\sqrt{\left\{ \frac{37028}{8} - (68)^2 \right\} \left\{ \frac{38132}{8} - (69)^2 \right\}}} \\ &= \frac{4695 - 4692}{\sqrt{(4628.5 - 4624)(4766.5 - 4761)}} \\ &= \frac{3}{\sqrt{4.5 \times 5.5}} = 0.603 \end{aligned}$$

Confidence Intervals

Under the assumption that the observations are normally and independent distributed, a $100(1-\alpha)\%$ confidence interval on the slope β_1 in simple linear regression is

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}$$

Similarly, a $100(1-\alpha)\%$ confidence interval on the intercept β_0 is

$$\hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]} \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}$$

Q We will find a 95% confidence interval on the slope of the regression line using data $\hat{\beta}_1 = 14.947$, $S_{xx} = 0.68088$ and $\hat{\sigma}^2 = 1.18$, $n = 20$.

Sm From eqn -

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}$$

$$14.947 - t_{0.025, 18} \sqrt{\frac{1.18}{0.68088}} \leq \beta_1 \leq 14.947 + t_{0.025, 18} \sqrt{\frac{1.18}{0.68088}}$$

$$14.947 - 2.101 \sqrt{\frac{1.18}{0.68088}} \leq \beta_1 \leq 14.947 + 2.101 \sqrt{\frac{1.18}{0.68088}}$$

This simplifies to

$$12.181 \leq \beta_1 \leq 17.713$$

Interpretation - The CI does not include zero, so there is strong evidence (at $\alpha = 0.05$) that the slope is not zero.

Hypothesis Tests in Simple Linear Regression -

Suppose we wish to test the hypothesis that the slope equals a constant, say $\beta_{1,0}$. The appropriate hypotheses are

$$H_0 : \beta_1 = \beta_{1,0}$$

$$H_1 : \beta_1 \neq \beta_{1,0}$$

Test statistic -

$$t_0 = \frac{\hat{\beta}_1 - \beta_{1,0}}{\sqrt{\hat{\sigma}^2 / S_{xx}}}$$

$$\text{where, } S_{xx} = \sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}$$

Follow the t distribution with $n-2$ degrees of freedom under $H_0 : \beta_1 = \beta_{1,0}$. We would reject $H_0 : \beta_1 = \beta_{1,0}$ if

$$|t_0| > t_{\alpha/2, n-2}$$

A similar procedure can be used to test hypotheses about the intercept.

$$H_0 : \beta_0 = \beta_{0,0}$$

$$H_1 : \beta_0 \neq \beta_{0,0}$$

Test statistic -

$$t_0 = \frac{\hat{\beta}_0 - \beta_{0,0}}{\sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}}$$

and reject the null hypothesis if the computed value of this test statistic, t_0 , is such that $|t_0| > t_{\alpha/2, n-2}$.

Q Using the estimated value $\hat{\beta}_1 = 0.903643$, test the hypothesis $\beta_1 = 1.0$ against the alternative that $\beta < 1.0$.

Soln

Null Hypothesis $\beta_1 = 1.0$

Alternative Hypothesis $\beta_1 < 1.0$

Given - $S_{xx} = 4152.18$

$\sigma = 3.2295$

$n = 33$

$t_{0.05, 31} = 1.648$

Test statistic $t = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}}$

$$t = \frac{0.903643 - 1.0}{3.2295 / \sqrt{4152.18}}$$

$$t = -1.92$$

with $n-2 = 31$ degree of freedom

Decision: Since $t > t_{\alpha, n-2}$. We Reject the null hypothesis.

Suggesting that it is a strong evidence that $\beta < 1.0$.