

Linear Regression:

①
- Dr. Ash Singh
USAR, GGSIPU

Regression analysis is statistical method used to estimate relationship between two or more random variables.

In regression model, there is one dependent random variable (response variable), and one or more independent variables (predictors).

r.v.s. | $Y \rightarrow$ Response variable
 $X^{(1)}, X^{(2)}, \dots, X^{(k)} \rightarrow$ predictors.

Linear regression model

$$Y = \beta_0 + \beta_1 X^{(1)} + \beta_2 X^{(2)} + \dots + \beta_k X^{(k)} + \epsilon$$

Here ϵ is a random error variable, with mean zero and unknown variance σ^2 . ↳ Error

Consider the case when only one predictor is there —

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

This is called Univariate Linear Regression model.

Suppose that we have n pairs of observations

(2)

$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$.

Then we have —

$$y_1 = \beta_0 + \beta_1 x_1 + \epsilon_1$$

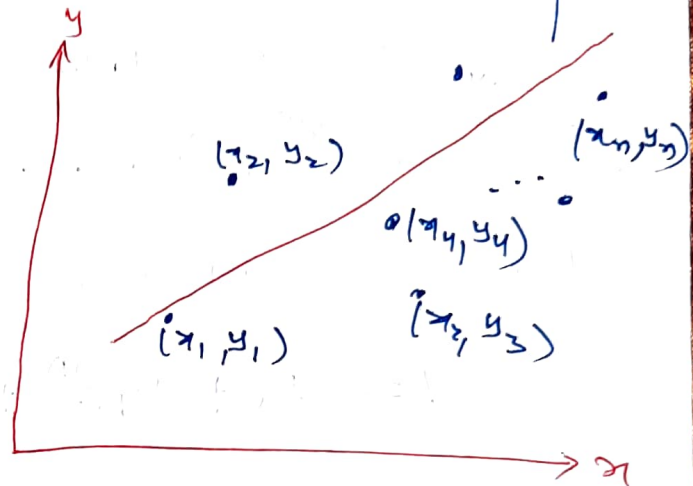
$$y_2 = \beta_0 + \beta_1 x_2 + \epsilon_2$$

\vdots

$$y_n = \beta_0 + \beta_1 x_n + \epsilon_n$$

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

↑
Line of Regression



Consider following sum of squares of error terms —

$$L = \sum_{i=1}^n \epsilon_i^2$$

→ [It is also sum of squares of the deviations of observations from the true regression line.]

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

(3)

Least square linear regression model:

To find the estimates of slope ($\hat{\beta}_1$) and intercept ($\hat{\beta}_0$) is the fitted line of regression ($\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$), we minimize function L with respect to β_0 & β_1 .

From theory of calculus, we get

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

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

Simplifying above two equations:

$$\sum_{i=1}^n y_i - \hat{\beta}_0 \sum_{i=1}^n 1 - \hat{\beta}_1 \sum_{i=1}^n x_i = 0 \quad \text{--- ①}$$

$$\sum_{i=1}^n x_i y_i - \hat{\beta}_0 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0 \quad \text{--- ②}$$

From ①

$$\sum y_i - \hat{\beta}_0 n - \hat{\beta}_1 \sum x_i = 0$$

$$\Rightarrow \hat{\beta}_0 = \frac{\sum y_i}{n} - \hat{\beta}_1 \frac{\sum x_i}{n}$$

$$\boxed{\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}} \quad \text{--- (A)}$$

From ②

④

$$\Rightarrow \sum x_i y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) \sum x_i - \hat{\beta}_1 \sum x_i^2 = 0$$

$$\sum x_i y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) n \bar{x} - \hat{\beta}_1 \sum x_i^2 = 0$$

$$\hat{\beta}_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}$$

Also we can write

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

So fitted line of regression is —

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

where

$$\beta_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\text{where } \bar{y} = \frac{\sum y_i}{n}$$

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

$$\bar{x} = \frac{\sum x_i}{n}$$

Eg Do eg 11.1 of ~~Reference~~ ^{Text} book.

Sample Correlation Coefficient (R):

(5)

Consider observations (x_i, y_i) , $i=1, 2, \dots, n$ corresponding to the r.v.s X & Y .

The sample correlation coefficient (R) is given as follows —

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

We can see that

$$\hat{\beta}_1 = \sqrt{\frac{\sum (y_i - \bar{y})^2}{\sum (x_i - \bar{x})^2}} \times R$$

$$\hat{\beta}_1 = \sqrt{\frac{\frac{1}{n-1} \sum (y_i - \bar{y})^2}{\frac{1}{n-1} \sum (x_i - \bar{x})^2}} R$$

$$\hat{\beta}_1 = \frac{s_y}{s_x} R$$

where s_y & s_x are sample standard deviations corresponding to r.v.s Y & X , respectively.

Multiple Linear Regression Model -

A regression model that contains ^⑥ more than one regressor variable is called a multiple linear regression model.

For eg. -

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + E$$

Logistic Regression -

A regression model in which response variable (Y) takes on only two possible values, 0 and 1, is called as logistic regression model.