

# Simple Linear Regression

---

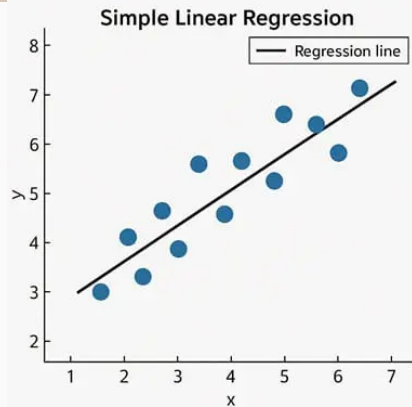
Dr. Mrityunjoy Barman, [mrityunjoybarman@soa.ac.in](mailto:mrityunjoybarman@soa.ac.in)

December 20, 2025

# Introduction to Simple Linear Regression i

- Simple linear regression is used to estimate the relationship between a **predictor variable** ( $x$ ) and a **response variable** ( $y$ ).
- It provides a linear approximation of how  $y$  changes as  $x$  changes.
- **Example:** Estimating the nutritional rating of cereals based on their sugar content.

## Introduction to Simple Linear Regression ii



**Figure 1:** Data fitting with a straight line.

# The Regression Equation

Let us consider the given data as:

$X$	$x_1$	$x_2$	$\cdots$	$x_n$
$Y$	$y_1$	$y_2$	$\cdots$	$y_n$

**Table 1:** The given data is given as  $(x_i, y_i)$ .

The estimated regression line is defined by

$$\hat{y} = b_0 + b_1x \quad (1)$$

- $\hat{y}$ : The estimated value of the response variable.
- $b_0$ : The y-intercept (estimated value of  $y$  when  $x = 0$ ).
- $b_1$ : The slope (estimated change in  $y$  per unit increase in  $x$ ).
- $b_0$  and  $b_1$  are called the **regression coefficients**.

## Method of Normal Equations

If the data points  $(x_i, y_i)$  were lying on the regression line (1), then we will have

$$y_i = b_0 + b_1 x_i$$

for all  $i = 1, 2, \dots, n$ .

This can be viewed as

$$\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \dots & \dots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix}, \quad \text{which gives } Pq = Q. \quad (2)$$

Multiplying by  $P^T$  bothsides we get,  $P^T Pq = P^T b$ . This equations are known as the normal equations. Now we can solve for the unknown vector  $q = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}$ .

# The Least-Squares Estimates i

The goal is to minimize the **Sum of Squared Errors (SSE)**:

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

By differentiating with respect to  $b_0$  and  $b_0$  and setting to zero, we derive:

$$\begin{aligned} \frac{\partial}{\partial b_0}(SSE) &= 0, & \frac{\partial}{\partial b_1}(SSE) &= 0 \\ \implies \sum y_i &= nb_0 + b_1 \sum x_i, & \text{and } \sum x_i y_i &= b_0 \sum x_i + b_1 \sum x_i^2. \end{aligned}$$

### Slope Estimate ( $b_1$ )

$$b_1 = \frac{\sum x_i y_i - \frac{1}{n}(\sum x_i)(\sum y_i)}{\sum x_i^2 - \frac{1}{n}(\sum x_i)^2} = \frac{\text{cov}(X, Y)}{\text{var}(X)}.$$

### Intercept Estimate ( $b_0$ )

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

## Example 1

### Example

Consider the following data:

$X$	-1	1	2
$Y$	1	1	3

Here, we see  $E(X) = 2/3$ ,  $E(X^2) = 2$ ,  $E(Y) = 5/3$ ,  $E(XY) = 2$ , and hence

$$\text{cov}(X, Y) = E(XY) - E(X)E(Y) = 2 - 10/9 = 8/9,$$

$$\text{and } \text{var}(X) = E(X^2) - E(X)^2 = 2 - 4/9 = 14/9,$$

$$\implies b_1 = \frac{\text{cov}(X, Y)}{\text{var}(X)} = \frac{8/9}{14/9} = 4/7.$$

Similarly,  $b_0 = \bar{y} - b_1\bar{x} = 5/3 - 4/7 \times 2/3 = 9/7$ .

Thus, the regression line is given by  $\hat{y} = 9/7 + 4/7x$ .



## Example 2: Calculation of the SSE

The SSE represents the overall measure of prediction error. Below is the calculation for 10 competitors using  $\hat{y} = 6 + 2x$ .

Subject	Time ( $x$ )	Distance ( $y$ )	Predicted ( $\hat{y}$ )	Residual ( $y - \hat{y}$ )	$(y - \hat{y})^2$
1	2	10	10	0	0
2	2	11	10	1	1
3	3	12	12	0	0
4	4	13	14	-1	1
5	4	14	14	0	0
6	5	15	16	-1	1
7	6	20	18	2	4
8	7	18	20	-2	4
9	8	22	22	0	0
10	9	25	24	1	1
50		160			SSE=12.

## Various types of Estimation Error: Measuring Goodness of Fit i

- **Sum of Squared Error**  $SSE = \sum (y - \hat{y})^2$ .
- **Sum of Squared Total**  $SST = \sum (y - \bar{y})^2$
- **Sum of Squared Regression**  $SSR = SST - SSE$ .
- We also call the constants  $b_0$ ,  $b_1$  as the regression coefficients.

The **Coefficient of Determination** ( $r^2$ ) measures the proportion of variability in  $y$  explained by the regression.

$$r^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}.$$

- $SST = SSR + SSE$ .

- $r^2$  ranges from 0 to 1.
- Values near 1 indicate an extremely good fit.

## Standard Error of the Estimate ( $s$ ): Correlation Coefficient ( $r$ )

The Mean Squared Error and the standard error  $s$  are defined as

$$MSE = \frac{SSE}{n-2}, \quad s = \sqrt{MSE} = \sqrt{\frac{SSE}{n-2}}.$$

The **Pearson correlation coefficient** ( $r$ ) measures the strength and direction of the linear relationship.

$$\rho = \frac{\sum xy - (\sum x)(\sum y)/n}{\sqrt{\sum x^2 - (\sum x)^2/n} \sqrt{\sum y^2 - (\sum y)^2/n}} = \frac{E(XY) - E(X)E(Y)}{\sigma_X \sigma_Y}$$

- Range:  $[-1, 1]$ .
- Positive  $\rho$ :  $y$  increases as  $x$  increases.
- Negative  $\rho$ :  $y$  decreases as  $x$  increases.
- $\rho = \pm\sqrt{r^2}$  (sign depends on the slope  $b_1$ ).
- We can also express it as  $\rho = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$ .

## Example 3: Projected Score

### Example

Suppose, in a T20 match between India and South Africa, the progress of runs scored in India innings are given as follows:

Over	4	8	12	16
Runs	33	68	115	150

- (a) Find the estimated linear regression line to the above data.
- (b) What is the projected score at the end of the India innings?
- (c) What are various types of errors in this estimation?
- (d) How good were the above data fit by the regression line? Explain using the coefficient of determination.
- (e) Find the correlation coefficient  $\rho$  between the attributes "**Over**" and "**Runs**".

## Example 4 i

We already found that the regression line for the following data is given by  $\hat{y} = 6 + 2x$ .  
Note that  $\bar{x} = 5$ ,  $\bar{y} = 16$ .

Subject	Time ( $x$ )	Distance ( $y$ )	Predicted( $\hat{y}$ )	$(y - \bar{y})$	$(y - \bar{y})^2$
1	2	10	10	-6	36
2	2	11	10	-5	25
3	3	12	12	-4	16
4	4	13	14	-3	9
5	4	14	14	-2	4
6	5	15	16	-1	1
7	6	20	18	4	16
8	7	18	20	2	4
9	8	22	22	6	36
10	9	25	24	9	81
50		160		SST=228	

## Example 4 ii

Subject	Time ( $x$ )	Distance ( $y$ )	Predicted ( $\hat{y}$ )	$(\hat{y} - \bar{y})$	$(\hat{y} - \bar{y})^2$
1	2	10	10	-6	36
2	2	11	10	-6	36
3	3	12	12	-4	16
4	4	13	14	-2	4
5	4	14	14	-2	4
6	5	15	16	0	0
7	6	20	18	2	4
8	7	18	20	4	16
9	8	22	22	6	36
10	9	25	24	8	64
50		160			SSR=216

Hence, we find the coefficient of determination  $r^2 = \frac{SSR}{SST} = \frac{216}{228} = 0.947$ .

## Example 5: Predicted Score in Chemistry

### Example

Compute and interpret the correlation coefficient for the following data:

Mathematics score %	70	92	80	74	65	83
Chemistry score (%)	74	84	63	87	78	90

- (a) Find the estimated linear regression line to the above data.
- (b) What is the predicted score in Chemistry if the student has 95% marks in Mathematics?
- (c) What are various types of errors in this estimation?
- (d) How good were the above data fit by the regression line?



## Example 6: Pressure-Volume relation in a gas

### Example

Suppose, the volume  $V$  of a non-ideal gas subject to various pressure  $P$  are recorded as follows:

$P$ in ( $kg/cm^2$ )	64.7	51.3	40.5	25.9	7.8
$V$ in ( $cm^3$ )	50	60	70	90	100

- (a) The ideal gas follows the certain law of type  $PV^\gamma = C$ , where  $\gamma$ ,  $C$  are constants. Estimate these constants  $\gamma$ ,  $C$ .
- (b) What is the predicted volume of the gas if we increase the pressure to  $90 \text{ kg/cm}^2$ .