

Linear Regression

1. Introduction

Linear regression is a statistical technique used to model the relationship between a dependent variable (response) and one or more independent variables (predictors). The key idea is to approximate this relationship by a linear function.

- **Dependent variable (Y):** also called outcome, response, or target.
- **Independent variable (X):** also called predictor, explanatory, or regressor.

Regression is used for:

- Prediction
 - Estimation
 - Hypothesis testing
 - Modeling causal relationships
-

2. The Simple Linear Regression Model

For one predictor variable:

$$y = \beta_0 + \beta_1 x + \epsilon$$

where:

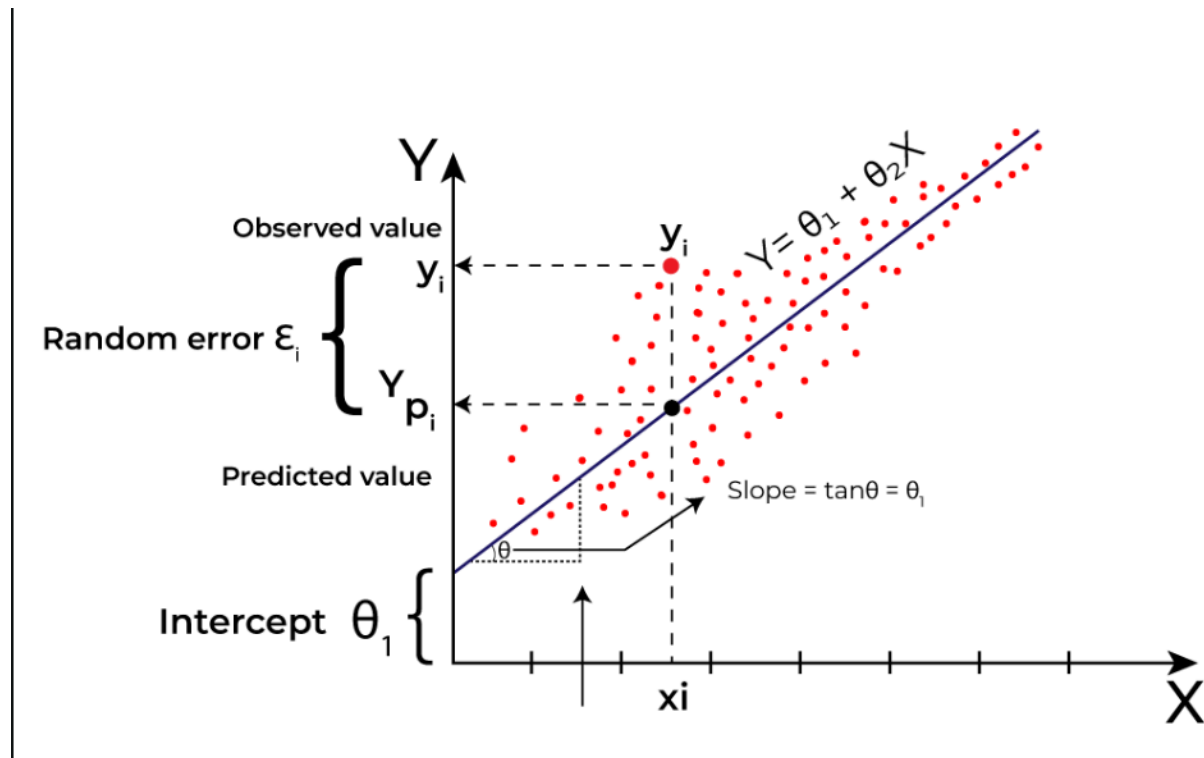
- Y is the dependent variable (the value we want to predict).
- X is the independent variable.
- β_0 is the **y-intercept**, the value of Y when $X=0$.
- β_1 is the **slope**, representing the change in Y for a one-unit change in X.
- ϵ is the error term, accounting for the variability in Y that is not explained by X.

Assumptions of Linear Regression

1. **Linearity:** Relationship between predictors and target is linear.
2. **Independence:** Observations are independent of each other.
3. **Homoscedasticity:** Constant variance of errors across values of predictors.

4. **Normality of Errors:** Errors (residuals) follow a normal distribution.
5. **No Multicollinearity:** Independent variables should not be highly correlated.

Example figure:



This is called a probabilistic model because for any fixed x , y is not exactly equal to $\beta_0 + \beta_1 x$ but varies around it.

Expected value:

$$E(Y | x) = \beta_0 + \beta_1 * x$$

Variance:

$$\text{Var}(Y | x) = \sigma^2$$

3. Example

Suppose the regression line is:

$$y = 7.5 + 0.5 * x$$

- Interpretation of slope $\beta_1 = 0.5$: On average, for every 1 unit increase in x, y increases by 0.5.
 - If $x = 20$, expected y is $7.5 + 0.5*20 = 17.5$.
-

4. Estimating Parameters (Least Squares)

We estimate β_0 and β_1 by minimizing the sum of squared errors:

$$SSE = \sum (y_i - (\beta_0 + \beta_1 x_i))^2$$

The estimates are:

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

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 * \bar{x}$$

5. Predicted and Residual Values

For each observation x_i , the predicted value is:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 * x_i$$

The residual is:

$$e_i = y_i - \hat{y}_i$$

Residuals are the difference between observed and predicted values.

6. Error Variance and R-squared

Error variance estimate:

$$\hat{\sigma}^2 = SSE / (n - 2)$$

where n is the number of observations.

Coefficient of determination:

$$R^2 = 1 - SSE / SST$$

where $SST = \sum (y_i - \bar{y})^2$.

R^2 measures the proportion of variance in y explained by the regression model.

7. Multiple Linear Regression

Extension of simple regression to multiple predictors:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$$

- Each β_i is a partial regression coefficient: the effect on y when x_i increases by 1, holding other predictors constant.