

Linear Regression in Machine Learning

Linear regression is a type of [supervised machine-learning algorithm](#) that learns from the labelled datasets and maps the data points with most optimized linear functions which can be used for prediction on new datasets. It assumes that there is a linear relationship between the input and output, meaning the output changes at a constant rate as the input changes. This relationship is represented by a straight line.

For example, we want to predict a student's exam score based on how many hours they studied. We observe that as students study more hours, their scores go up. In the example of predicting exam scores based on hours studied.

Independent variable (input): Hours studied because it's the factor we control or observe.

- Dependent variable (output): Exam score because it depends on how many hours were studied.

We use the independent variable to predict the dependent variable.

What is Linear Regression?

Linear regression is a statistical method used to model the relationship between a dependent variable and one or more independent variables.

Formula : $\hat{y} = \theta_0 + \theta_1 x$

\hat{y} : The predicted value of the dependent variable (output).

x : Input variable

Example:

Predicting a person's salary based on years of experience.



How Linear Regression Works

$$\min \frac{1}{n} \sum_{i=1}^n (\bar{y}_i - y_i)^2$$

Where:

\bar{y}_i : Predicted value for the i^{th} observation

y_i : Actual observed value for the i^{th} observation

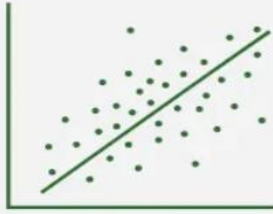
n : Total number of observations

Min: Indicates the goal is to minimize the total error

The model adjusts the parameters $\theta_0, \theta_1, \dots, \theta_n$ to minimize this error

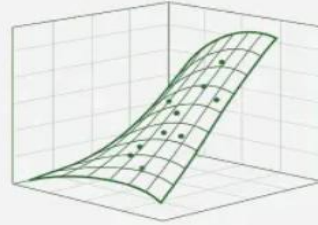
Types of Linear Regression

Simple Linear Regression



Predicts the dependent variable using a single independent variable.

Multiple Linear Regression



Uses two or more independent variables to predict the dependent variable.

Real-World Use Cases of Linear Regression

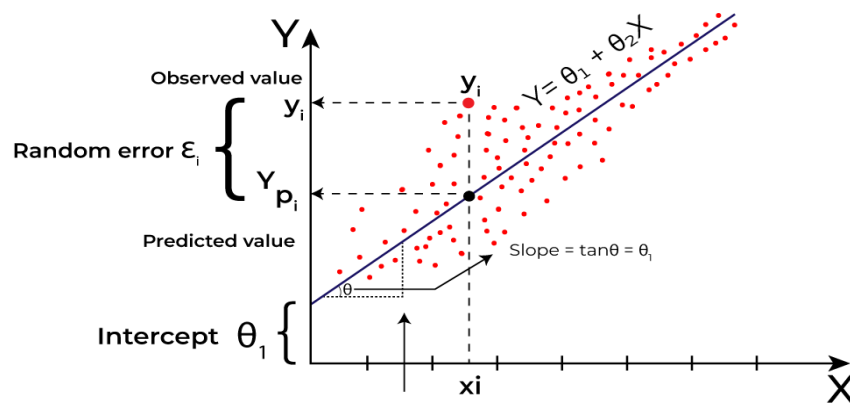


Best Fit Line in Linear Regression

In linear regression, the best-fit line is the straight line that most accurately represents the relationship between the independent variable (input) and the dependent variable (output). It is the line that minimizes the difference between the actual data points and the predicted values from the model.

1. Goal of the Best-Fit Line

The goal of linear regression is to find a straight line that minimizes the error (the difference) between the observed data points and the predicted values. This line helps us predict the dependent variable for new, unseen data.



Linear Regression

Here Y is called a dependent or target variable and X is called an independent variable also known as the predictor of Y. There are many types of functions or modules that can be used for regression. A linear function is the simplest type of function. Here, X may be a single feature or multiple features representing the problem.

2. Equation of the Best-Fit Line

For simple linear regression (with one independent variable), the best-fit line is represented by the equation.

$$y = mx + b$$

Where:

- **y** is the predicted value (dependent variable)
- **x** is the input (independent variable)
- **m** is the slope of the line (how much y changes when x changes)
- **b** is the intercept (the value of y when x = 0)

The best-fit line will be the one that optimizes the values of m (slope) and b (intercept) so that the predicted y values are as close as possible to the actual data points.

3. Minimizing the Error: The Least Squares Method

To find the best-fit line, we use a method called [Least Squares](#). The idea behind this method is to minimize the sum of squared differences between the actual values (data points) and the predicted values from the line. These differences are called residuals.

The formula for residuals is:

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

Where:

- y_i is the actual observed value
- \hat{y}_i is the predicted value from the line for that x_i

The least squares method minimizes the sum of the squared residuals:

$$Sum of squared errors (SSE) = \sum (y_i - \hat{y}_i)^2$$

This method ensures that the line best represents the data where the sum of the squared differences between the predicted values and actual values is as small as possible.

4. Interpretation of the Best-Fit Line

- **Slope (m):** The slope of the best-fit line indicates how much the dependent variable (y) changes with each unit change in the independent variable (x).

For example, if the slope is 5, it means that for every 1-unit increase in x, the value of y increases by 5 units.
- **Intercept (b):** The intercept represents the predicted value of y when x = 0. It's the point where the line crosses the y-axis.

In linear regression some hypotheses are made to ensure reliability of the model's results.

Limitations:

- **Assumes Linearity:** The method assumes the relationship between the variables is linear. If the relationship is non-linear, linear regression might not work well.
- **Sensitivity to Outliers:** Outliers can significantly affect the slope and intercept, skewing the best-fit line.

Hypothesis function in Linear Regression

In linear regression, the hypothesis function is the equation used to make predictions about the dependent variable based on the independent variables. It represents the relationship between the input features and the target output.

For a simple case with one independent variable, the hypothesis function is:

$$h(x) = \beta_0 + \beta_1 x$$

Where:

- $h(x)$ (or \hat{y}) is the predicted value of the dependent variable (y).
- x is the independent variable.
- β_0 is the intercept, representing the value of y when x is 0.
- β_1 is the slope, indicating how much y changes for each unit change in x .

For **multiple linear regression** (with more than one independent variable), the hypothesis function expands to:

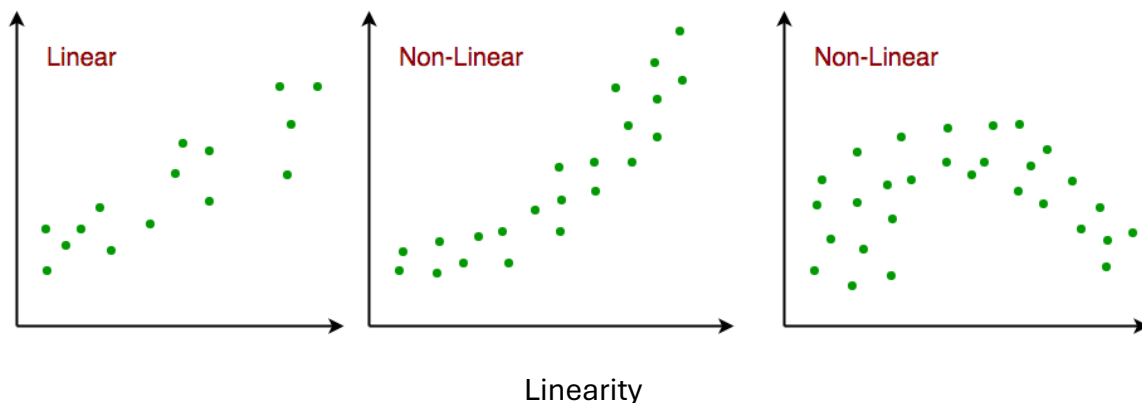
$$h(x_1, x_2, \dots, x_k) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

Where:

- x_1, x_2, \dots, x_k are the independent variables.
- β_0 is the intercept.
- $\beta_1, \beta_2, \dots, \beta_k$ are the coefficients, representing the influence of each respective independent variable on the predicted output.

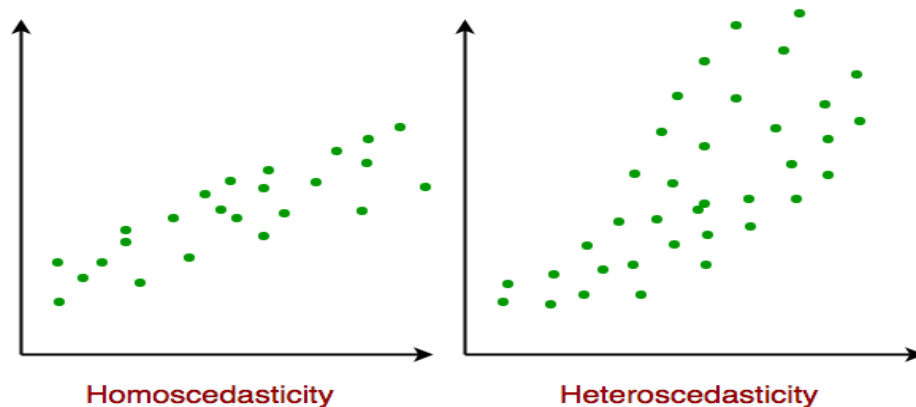
Assumptions of the Linear Regression

1. Linearity: The relationship between inputs (X) and the output (Y) is a straight line.



2. Independence of Errors: The errors in predictions should not affect each other.

3. Constant Variance (Homoscedasticity): The errors should have equal spread across all values of the input. If the spread changes (like fans out or shrinks), it's called heteroscedasticity and it's a problem for the model.



Homoscedasticity

4. Normality of Errors: The errors should follow a normal (bell-shaped) distribution.

5. No Multicollinearity (for multiple regression): Input variables shouldn't be too closely related to each other.

6. No Autocorrelation: Errors shouldn't show repeating patterns, especially in time-based data.

7. Additivity: The total effect on Y is just the sum of effects from each X, no mixing or interaction between them.

Cost function for Linear Regression

As we have discussed earlier about best fit line in linear regression, it's not easy to get it easily in real life cases so we need to calculate errors that affects it. These errors need to be calculated to mitigate them. The difference between the predicted value \hat{Y} and the true value Y and it is called [cost function](#) or the [loss function](#).

In Linear Regression, the Mean Squared Error (MSE) cost function is employed, which calculates the average of the squared errors between the predicted values \hat{y}_i and the actual values y_i . The purpose is to determine the optimal values for the intercept θ_1 and the coefficient of the input feature θ_2 providing the best-fit line for the given data points. The linear equation expressing this relationship is

$$\hat{y}_i = \theta_1 + \theta_2 x_i$$

MSE function can be calculated as:

$$\text{Cost function}(J) = \frac{1}{n} \sum_n^i (\hat{y}_i - y_i)^2$$

Now we have calculated loss function we need to optimize model to mitigate this error and it is done through gradient descent.