

Machine learning - Linear regression

Botond Koroknai

2025.08.24

1 Introduction

Linear regression is one of the most fundamental and widely used models in the field of statistics and machine learning.

It can be used to model the relationship between input variables (features) and a continuous output variable (target).

At first, we have to assume that this relationship can be approximated by a straight line (if we are talking about a single feature) or a hyperplane (in case of multiple features).

1.1 Why is linear regression considered fundamental in the world of machine learning?

- **Simplicity:** The mathematical formulation of linear regression is straightforward, making it an ideal starting point for understanding more complex models.
- **Interpretability:** The learned parameters (weights) directly indicate the contribution of each input feature to the output prediction, providing clear insights into the data.
- **Practical relevance:** Despite its simplicity, linear regression is highly effective for many real-world problems where relationships are approximately linear or where interpretability is essential.

2 Mathematical background

As previously mentioned, the idea behind using linear regression is to approximate the relationship between the inputs (features) and the output (target) using a linear function.

This way the output can be expressed as a weighted combination of the input features plus a constant, which we call bias.

2.1 Single Feature (1D case):

$$\hat{y} = wx + b$$

where

- w is the slope of the line, determining how strongly x influences the output,
- b represents the value of \hat{y} when $x = 0$,
- \hat{y} is the predicted value for the given input x .

In this case, the model corresponds to finding the best line that passes as closely as possible through the data points in a two-dimensional plane.

2.2 Multiple Features (vector form)

When we have more than one input feature

$$x = (x_1, x_2, \dots, x_d)$$

the linear model generalizes to:

$$\hat{y} = x_1w_1 + x_2w_2 + \dots + x_dw_d + b$$

Which can be written up in a vector notation as:

$$\hat{y} = \mathbf{x}^\top \mathbf{w} + b$$

where

- $\mathbf{x} \in \mathbb{R}^d$ is the feature vector,
- $\mathbf{w} \in \mathbb{R}^d$ is the weight vector,
- $b \in \mathbb{R}$ is the bias term.

Finally, when we consider all n training examples together, it is convenient to organize them into a single data matrix

$$X \in \mathbb{R}^{n \times d}$$

where each row corresponds to one training example and each column corresponds to a particular feature across all examples:

$$X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(n)} & x_2^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times d}$$

The weight vector $\mathbf{w} \in \mathbb{R}^d$ contains one parameter for each feature, and $b \in \mathbb{R}$ is the bias term. Using matrix multiplication, we can compute the predicted outputs for all training examples simultaneously as

$$\hat{\mathbf{y}} = X\mathbf{w} + b\mathbf{1}$$

where $\hat{\mathbf{y}} \in \mathbb{R}^n$ is the vector of predictions.

Using such a representation comes with several advantages:

- **Efficiency:** Modern numerical libraries (like NumPy) are optimized for matrix operations, allowing all predictions to be computed at once without explicit loops.
- **Clarity:** The notation clearly expresses the relationship between the data, the weights, and the predictions in a single formula.

In any case we are interested in a specific prediction $\hat{y}^{(i)}$ for a single example i , it can be quickly calculated using the

$$\hat{y}^{(i)} = (\mathbf{x}^{(i)})^\top \mathbf{w} + b$$

formula.

3 Loss Function and Cost Function

In order to train a linear regression model, we must find out how accurately the model predicts the data. This process consists of two steps: first, defining a loss function for a single example. Then, defining a cost function for the entire dataset.

3.1 Loss Function

The **loss function** quantifies the error of the model on a *single training example*. For linear regression, the most common choice is the **squared error**:

$$\mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = (\hat{y}^{(i)} - y^{(i)})^2$$

where $y^{(i)}$ is the true target value for the i -th example, and $\hat{y}^{(i)}$ is the corresponding prediction.

3.2 Cost Function

The **cost function** aggregates the loss over the entire training dataset. It is essentially the average of all individual losses, providing a single scalar value that represents the model's performance across all examples:

$$J(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})^2$$

3.3 Training Objective

The training objective is to find the parameters \mathbf{w} and b that *minimize the cost function*:

$$(\mathbf{w}^*, b^*) = \arg \min_{\mathbf{w}, b} \underbrace{\frac{1}{n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})^2}_{\text{Cost function } J(\mathbf{w}, b)}$$

Minimizing the cost function ensures that the model's predictions are as close as possible to the true target values on average across all examples.

4 Numerical Interpretation

Once we have defined the loss function and the training objective, the next step is to determine how to find the optimal parameters (\mathbf{w}^*, b^*) that minimize the cost function. There are two standard approaches in linear regression: a closed-form solution and an iterative approach.

4.1 Closed-Form Solution (Normal Equation)

For linear regression with the squared error loss, we can compute the exact optimal parameters directly using linear algebra. This method is called the **normal equation**.

Derivation

The cost function is

$$J(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})^2 = \frac{1}{n} \sum_{i=1}^n ((\mathbf{x}^{(i)})^\top \mathbf{w} + b - y^{(i)})^2$$

To find the minimum, we take the partial derivatives of J with respect to \mathbf{w} and b and set them equal to zero:

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{2}{n} \sum_{i=1}^n \mathbf{x}^{(i)} ((\mathbf{x}^{(i)})^\top \mathbf{w} + b - y^{(i)}) = 0$$

$$\frac{\partial J}{\partial b} = \frac{2}{n} \sum_{i=1}^n ((\mathbf{x}^{(i)})^\top \mathbf{w} + b - y^{(i)}) = 0$$

The second equation gives the bias directly:

$$b^* = \bar{y} - \bar{\mathbf{x}}^\top \mathbf{w}^*, \quad \text{where } \bar{y} = \frac{1}{n} \sum_{i=1}^n y^{(i)}, \quad \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}.$$

Substituting b^* into the first derivative equation, and using matrix notation, we obtain:

$$X^\top (X\mathbf{w}^* + b^*\mathbf{1} - \mathbf{y}) = 0$$

Centering the data by subtracting the mean feature vector $\bar{\mathbf{x}}$ and the mean target \bar{y} simplifies the equation to:

$$\tilde{X}^\top \tilde{X} \mathbf{w}^* = \tilde{X}^\top \tilde{\mathbf{y}}$$

where $\tilde{X} = X - \mathbf{1}\bar{\mathbf{x}}^\top$ and $\tilde{\mathbf{y}} = \mathbf{y} - \bar{y}\mathbf{1}$. Finally, solving for \mathbf{w}^* gives the closed-form solution:

$$\mathbf{w}^* = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{\mathbf{y}}$$

and the bias is recovered as

$$b^* = \bar{y} - \bar{\mathbf{x}}^\top \mathbf{w}^*$$

This derivation shows that the optimal weights are found by projecting the target vector \mathbf{y} onto the space formed by the features, which gives predictions that are as close as possible to the true values in terms of squared error.

4.2 Iterative Solution (Gradient Descent)

Gradient descent is an iterative method to find the optimal parameters (\mathbf{w}^*, b^*) that minimize the cost function. Instead of solving a closed-form equation, we start with an initial guess for \mathbf{w} and b , and iteratively adjust them to reduce the prediction error. As a reminder, the cost function is

$$J(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n ((\mathbf{x}^{(i)})^\top \mathbf{w} + b - y^{(i)})^2$$

4.2.1 Step 1: Compute gradients

To update the parameters, we compute the partial derivatives of the cost function with respect to \mathbf{w} and b . Using calculus:

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left[\frac{1}{n} \sum_{i=1}^n ((\mathbf{x}^{(i)})^\top \mathbf{w} + b - y^{(i)})^2 \right] = \frac{2}{n} \sum_{i=1}^n ((\mathbf{x}^{(i)})^\top \mathbf{w} + b - y^{(i)}) \mathbf{x}^{(i)} \\ \frac{\partial J}{\partial b} &= \frac{\partial}{\partial b} \left[\frac{1}{n} \sum_{i=1}^n ((\mathbf{x}^{(i)})^\top \mathbf{w} + b - y^{(i)})^2 \right] = \frac{2}{n} \sum_{i=1}^n ((\mathbf{x}^{(i)})^\top \mathbf{w} + b - y^{(i)}) \end{aligned}$$

These gradients quantify how much the cost function changes with respect to each parameter, also they indicate the direction in which the parameters should be adjusted to reduce the error.

4.2.2 Step 2: Update rules

We update the parameters by moving in the direction opposite to the gradient:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial J}{\partial \mathbf{w}}, \quad b \leftarrow b - \eta \frac{\partial J}{\partial b}$$

where $\eta > 0$ is the learning rate, controlling the size of each step.

4.2.3 Step 3: Matrix representation

For computational efficiency, the gradient can be expressed in matrix notation. Let $\hat{\mathbf{y}} = X\mathbf{w} + b\mathbf{1}$ be the vector of predictions for all n training examples. Then the gradients become:

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{2}{n} X^\top (\hat{\mathbf{y}} - \mathbf{y}), \quad \frac{\partial J}{\partial b} = \frac{2}{n} \mathbf{1}^\top (\hat{\mathbf{y}} - \mathbf{y})$$

where $X \in \mathbb{R}^{n \times d}$ is the feature matrix, $\mathbf{y} \in \mathbb{R}^n$ is the vector of target values, and $\mathbf{1} \in \mathbb{R}^n$ is a vector of ones.

By substituting back the values, the update rules in matrix form are then:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{2}{n} X^\top (\hat{\mathbf{y}} - \mathbf{y}), \quad b \leftarrow b - \eta \frac{2}{n} \mathbf{1}^\top (\hat{\mathbf{y}} - \mathbf{y})$$

The overall goal is to move iteratively in the direction opposite to the gradient, thereby minimising the cost function.