

Linear Regression

Epoch IIT Hyderabad

Ankit Saha
AI21BTECH11004

15 Jul 2022

1. INTRODUCTION

Linear regression is a machine learning algorithm that predicts real-valued labels for a given input. Estimating real-valued output is known as regression and estimating categorical output is known as classification. Linear regression is one of the easiest and most-commonly used regression algorithms.

Linear regression tries to find a linear relationship between a dependent variable and one or more independent variables (also known as regressors).

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_n x_n \quad (1.1)$$

This can be written in matrix notation as

$$y = \boldsymbol{\beta}^\top \mathbf{x} \quad (1.2)$$

where $\boldsymbol{\beta} = (\beta_0 \ \beta_1 \ \cdots \ \beta_n)^\top$ and $\mathbf{x} = (1 \ x_1 \ \cdots \ x_n)^\top$

Thus, the objective of linear regression is to find the optimal parameter vector $\boldsymbol{\beta}$, i.e. finding the line (or hyperplane) that best fits the given data.

2. MATHEMATICAL FORMULATION

Suppose we are given m labelled samples $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)$ where $x_i \in \mathbb{R}^{n+1}, y_i \in \mathbb{R} \ \forall i \in [m]$. The predicted output matrix is given by $\hat{\mathbf{y}}^\top = \boldsymbol{\beta}^\top \mathbf{X}$ or $\hat{\mathbf{y}} = \mathbf{X}^\top \boldsymbol{\beta}$ where

$$\hat{\mathbf{y}} = (\hat{y}_1 \ \hat{y}_2 \ \cdots \ \hat{y}_m)^\top \quad (2.1)$$

$$\boldsymbol{\beta} = (\beta_0 \ \beta_1 \ \cdots \ \beta_n)^\top \quad (2.2)$$

$$\mathbf{X} = (\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_m) = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_{11} & x_{21} & \cdots & x_{m1} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1n} & x_{2n} & \cdots & x_{mn} \end{pmatrix} \quad (2.3)$$

The average least squares error is often used as an evaluation metric for linear regression. It is given by

$$J(\boldsymbol{\beta}) = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \quad (2.4)$$

$$= \frac{1}{m} \|\hat{\mathbf{y}} - \mathbf{y}\|^2 \quad (2.5)$$

$$= \frac{1}{m} \|\mathbf{X}^\top \boldsymbol{\beta} - \mathbf{y}\|^2 \quad (2.6)$$

$$= \frac{1}{m} (\mathbf{X}^\top \boldsymbol{\beta} - \mathbf{y})^\top (\mathbf{X}^\top \boldsymbol{\beta} - \mathbf{y}) \quad (2.7)$$

Our goal is to minimize this average least squares error. The optimal parameter vector β is given by

$$\hat{\beta} = \arg \min_{\beta} J(\beta) \quad (2.8)$$

$$= \arg \min_{\beta} \frac{1}{m} (\beta^T \mathbf{X} - \mathbf{y}^T) (\mathbf{X}^T \beta - \mathbf{y}) \quad (2.9)$$

$$= \arg \min_{\beta} \frac{1}{m} (\beta^T \mathbf{X} \mathbf{X}^T \beta - \beta^T \mathbf{X} \mathbf{y} - \mathbf{y}^T \mathbf{X}^T \beta + \mathbf{y}^T \mathbf{y}) \quad (2.10)$$

3. FINDING THE OPTIMAL PARAMETERS

On differentiating $J(\beta)$ with respect to β (the terms are all scalars),

$$\frac{\partial}{\partial \beta} (\beta^T \mathbf{X} \mathbf{X}^T \beta) = \beta^T ((\mathbf{X} \mathbf{X}^T)^T + (\mathbf{X} \mathbf{X}^T)) = 2\beta^T \mathbf{X} \mathbf{X}^T = 2\mathbf{X} \mathbf{X}^T \beta \quad (3.1)$$

$$\frac{\partial}{\partial \beta} (\beta^T \mathbf{X} \mathbf{y}) = \mathbf{X} \mathbf{y} \quad (3.2)$$

$$\frac{\partial}{\partial \beta} (\mathbf{y}^T \mathbf{X}^T \beta) = \mathbf{X} \mathbf{y} \quad (3.3)$$

$$\frac{\partial}{\partial \beta} (\mathbf{y}^T \mathbf{y}) = \mathbf{0} \quad (3.4)$$

$$\frac{\partial}{\partial \beta} (J(\beta)) = 2\mathbf{X} \mathbf{X}^T \beta - 2\mathbf{X} \mathbf{y} \quad (3.5)$$

On equating the derivative to zero, we get

$$\mathbf{X} \mathbf{X}^T \hat{\beta} = \mathbf{X} \mathbf{y} \quad (3.6)$$

$$\therefore \hat{\beta} = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{y} \quad (3.7)$$

4. COEFFICIENT OF DETERMINATION

In some cases, another evaluation metric known as the coefficient of determination (R^2 value) is used.

$$R^2 = 1 - \frac{\sum_{i=1}^m (y_i - \hat{y}_i)^2}{\sum_{i=1}^m (y_i - \bar{y})^2} \quad (4.1)$$

where \bar{y} is the mean of all $y_i, i \in [m]$

$$\bar{y} = \frac{1}{m} \sum_{i=1}^m y_i \quad (4.2)$$

The coefficient of determination tells us about how much of the variation in the dependent variable is explained by the variation in the independent variables. A high R^2 implies that the dependent variable is heavily dependent on the independent variables and our model is accurate and on the other hand, a low R^2 implies that the dependent variable is nearly independent of the dependent variables and our model is thus inaccurate.

5. QUESTIONS

- i) What does the parameter β_0 represent and why is it important?
- ii) Why can't we consider the sum of the differences instead of the sum of the square of the differences in the least squares method?
- iii) Similarly, why do we not use the sum of the absolute values of the differences?
- iv) What is the R^2 value when the number of samples, $m = 2$?
- v) Can linear regression be used for classification? Explain.

6. ANSWERS

- i) β_0 represents the y-intercept of the optimal line (or hyperplane). It is important because it provides the model another degree of freedom. If we do not consider β_0 , then the optimal line is constrained to pass through the origin, which may not be ideal. β_0 gives the model the freedom to start at any point it wants.
- ii) Differences can be both positive and negative. If we simply consider the differences, the positive differences may cancel out the negative differences. Thus, this does not give an accurate idea of the error of the model.
- iii) It is easier to solve for β analytically when we are using least squares error. If we consider absolute values, not only is it harder to solve because it is not differentiable, it also doesn't necessarily have a unique solution.
- iv) $m = 2$ means we are given 2 points as our data set. The best fitting line will be the line passing through these 2 points as we can always find a line passing through any 2 points. Thus, $y_i = \hat{y}_i \forall i$ and $R^2 = 1 - 0 = 1$
- v) Theoretically, linear regression can be used for classification by using a decision rule. But it is generally not used because linear regression is very sensitive to new data points. We will have to change our decision rule for every new data point added to keep the model accurate, which is not feasible in most cases.