

LINEAR REGRESSION

Nehal Khosla, Priyanshu Parida

National Institute of Science Education and Research

January 30, 2023

PART I: LINEAR REGRESSION: PROBLEM & SOLUTION

1	Regression	5
2	Linear Regression: The Problem	6
2.1	Introduction	6
2.2	Regression Line	7
2.3	Types	8
2.4	Mathematical Representation	9
2.5	Loss Function: Mean Squared Error	10
3	Linear Regression: The Solution	11
3.1	Solution	11
3.2	Quality of Fit	12

PART II: LINEAR REGRESSION: THE INDUCTIVE BIAS

- 1 Inductive Bias 14**
 - 1.1 A List 14
 - 1.2 Choice of Loss Function 15

PART III: LINEAR REGRESSION: APPLICATIONS AND SHORTCOMINGS

1 Applications and Shortcomings 17

Part I

LINEAR REGRESSION: THE PROBLEM

REGRESSION

- ▶ Regression is a statistical method that attempts to predict the strength and nature of the relationship between a dependent variable, and one or a series of independent variables.
- ▶ It does so by finding a curve that minimizes the error in the actual and expected values of the dependent variable of training data-set.
- ▶ For proper interpretation of regression, several assumptions (inductive bias) about the data and the model must hold.
- ▶ Linear regression is one of the common forms of this method. It establishes a linear relationship between the two variables.

LINEAR REGRESSION

INTRODUCTION

- ▶ Linear regression is a **supervised** machine learning algorithm.
- ▶ The model tries to find a **best-fit line** to establish a linear relationship between the dependent (y) and independent(x) variable.
- ▶ The model then uses this fit to predict the appropriate y-values for unknown x-values.
- ▶ The **best-fit line** is achieved minimizing the error between predicted and actual values. This is done by minimizing the **loss function**.

LINEAR REGRESSION

REGRESSION LINE

The line showing the linear relationship between the dependent and independent variable is called a **regression line**. An example of a regression line is shown below:

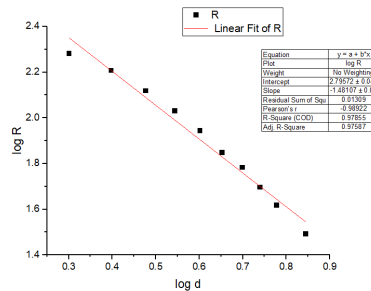


Figure. Regression line for log R (dependent variable) vs log d (independent variable)

The regression line may be positive, wherein the dependent variable increases with increase in independent variable; or it may be negative, wherein the dependent variable decreases with increase in independent variable (as in above figure).

LINEAR REGRESSION: THE PROBLEM

TYPES

Linear regression may be classified further into the following two types:

- ▶ **Simple Linear Regression:** Assumes a linear relationship between a single independent variable and a dependent variable.
- ▶ **Multiple Linear Regression:** Assumes a linear relationship between two or more independent variables and a dependent variable.

LINEAR REGRESSION: THE PROBLEM

MATHEMATICAL REPRESENTATION

Once a linear relationship has been determined by the algorithm, the general form of each model may be represented as follows:

▶ **Simple Linear Regression**

$$y = ax + b + u$$

▶ **Multiple Linear Regression**

$$Y = a_1x + a_2x + a_3x + b + u$$

where:

y = Dependent variable

x = Independent variable

a = Slope(s) of the variable(s)

b = The y-intercept

u = The regression residual/error term

LINEAR REGRESSION: THE PROBLEM

LOSS FUNCTION: MEAN SQUARED ERROR

- ▶ The regression line is achieved by minimizing the sum of mean squared error (loss function) for all points in the domain. The loss function is given as:

$$MSE = \frac{1}{N} \sum (y - f(x))^2$$

where $f(x) = a_1x + a_2x \dots + b$

LINEAR REGRESSION: THE SOLUTION

SOLUTION

The best fit line may be found in the following two manners:

► **Closed form (Exact form) Solution:**

- It solves the problem in terms of simple functions and mathematical operators.
- The closed form solution for linear regression is as follows:

$$B = (X'X)^{-1}X'Y$$

where B = Matrix of regression parameters

X = Matrix of X values

X' = Transpose of X

Y = Matrix of Y values

- Although this method gives an accurate model, it is computationally expensive, especially when there are more than 4 dimensions.

► **Gradient Descent:**

- It is used to minimize MSE by calculating the gradient of the loss function.
- It is an iterative optimization algorithm.

LINEAR REGRESSION: THE SOLUTION

QUALITY OF FIT

- ▶ The goodness of the fit achieved determines how linearly the variables are correlated.
- ▶ The goodness of fit may be calculated using the Pearson correlation coefficient, which is given by:

$$r = \frac{\Sigma(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\Sigma(x_i - \bar{x})^2 \Sigma(y_i - \bar{y})^2}}$$

- ▶ The higher the value of r , the better is the fit.

Part II

LINEAR REGRESSION: THE INDUCTIVE BIAS

INDUCTIVE BIAS

A LIST

Linear regression takes the following assumptions, or inductive biases:

- ▶ The assumption that the dependent and independent variables are linearly related.
- ▶ Homoscedasticity: The assumption that the error term should be the same for all points.
- ▶ The assumption that MSE is the most appropriate loss function for linear regression.

CHOICE OF LOSS FUNCTION

Let us analyse some loss functions to justify the choice of MSE as an appropriate loss function.

- ▶ $L_1 = (\mathbf{y}-\mathbf{f}(\mathbf{x}))$: This loss function gives out both positive and negative values, which cancel out to give near zero error for large data.
- ▶ $L_2 = |(\mathbf{y}-\mathbf{f}(\mathbf{x}))|$: Although errors do not cancel out here, the outliers are penalised equally as compared to standard data.
- ▶ $L_3 = (\mathbf{y}-\mathbf{f}(\mathbf{x}))^2$: In this case, the errors do not cancel out. Also, outliers are penalised more, giving a more appropriate regression line.

Hence, MSE is an appropriate choice for loss function.

Part III

LINEAR REGRESSION: APPLICATIONS AND SHORTCOMINGS

APPLICATIONS AND SHORTCOMINGS

- ▶ Linear regression finds its applications in several fields, like market analysis, financial analysis, environmental health, and medicine.
- ▶ However, it does leave some things to desire for. A linear correlation does not indicate causation, i.e. a connection between two variables does not imply that one causes the other.
- ▶ Linear regression is prone to noise and overfitting.
- ▶ It is prone to multicollinearity, i.e. occurrence of correlation between two or more independent variables. This reduces the statistical significance of an independent variable.

REFERENCES

1. CS460 Machine Learning 2023 Lectures, Subhankar Mishra.
2. Linear Regression in Machine Learning, Javatpoint.
3. ML | Linear Regression, geeksforgeeks.