

# Linear Regression: Deep Learning

Olaniyi Bayonle Alao

Summer Term, 2021

Bachelor of Electronic Engineering

Hochschule Hamm-Lippstadt

Lippstadt, Germany

olaniyi-bayonle.alao@stud.hshl.de

**Abstract**—talk about the data example prediction that will/was done using scikit-learn. Then talk about what linear regression, why it is used.....

this paper talks about linear regression in the context of deep learning. Linear regression is a statistical term that uses a dependent and independent variable to make predictions. in these paper, the use case of linear regression is used to predict data in machine learning. The data were gotten from kaggle.com. the data is a data of the titanic survivors and the model was trained using scikit-learn or tensorflow.

data cleaning was done using pandas a python library to prepare the data for model training.

regression analysis as a whole is a statistical method used to understand the relation between dependent and independent variables.

with the rise in the amount of data we have access to and increase in the performance of computers, Machine Learning a subset of Artificial intelligence has seen a significant growth. machine learning refers to the ability of applications to get better at doing things without necessarily a change in the code base then it is something that has really been helpful.

even though machine learning isn't a new concept, it is fast gaining recognition and changing lives thanks to the increase in the processing power of computers over the years to be able to process big data at a level that has never been experienced.

**Index Terms**—machine learning, linear regression, deep learning, scikit-learn

Unless otherwise stated, the main reference for facts, mathematical models and terminologies used in this research paper is gotten from [?]

## I. DIAGRAMS

Include diagram for machine learning automatically adapting to change?? get idea from fig 1-3 page 25 of [hands-on machine learning with scikit-learn]

## II. INTRODUCTION

Due to an explosive increase in the amount of data generated from different internet connected devices -Cyber-Physical-Systems (CPS)-, there has been an increase in the need to make sense of this "big data" to ensure the proper and productive functioning of businesses. By 2026, the data generation from global air fleet is projected to reach 98 billion gigabytes [?]. Thanks to the recent advancements in the information processing capabilities of computers, making sense of this huge chunk of data is now possible using data analysis techniques and Machine learning to clean and make prediction from this datasets. Machine learning is the ability of machine

-software application- to learn and make predictions without being programmed to do so. They are able to learn and get better with an increase in the amount of data they are being fed through their grounded mathematical foundations.

Linear regression is a statistical test used to find out the relationship between independent and dependent variables in a data set using mathematical formulae. They can also be used in projecting new relationship between the dependent and independent variables that has not been discovered. Even though linear models - linear regression - are quite simple to develop and understand, as well as good at predicting linear relationships, their approximation of nonlinear relationships have been found to be mostly unsatisfactory [?]. This paper goes more into the mathematical details this paper gives a profound, yet understandable mathematical description of Linear regression, relevant parts of the framework -Sci-Kit- used to perform predictions on the example dataset in the context of Machine learning using the Python programming language.

## III. THEORETICAL BACKGROUND

### A. Machine Learning

According to [?], "Machine learning is a branch of artificial intelligence (AI) focused on building applications that learn from data and improve their accuracy over time without being programmed to do so". Examples of applications built using machine learning trained models are email filters with the ability to distinguish between desired emails and spams, auto-suggestion/auto-correct in many typing applications, self-driving vehicles, amongst others. Machine learning techniques have the advantage of automatically adapting to change detected in trends from data. Likewise, they help reduce the complexity in writing application that are something's difficulty or even impossible to implement using traditional algorithms. Methods used in training a Machine learning model can be grouped into three main categories.

- Supervised Learning

This is a type of machine learning technique in which the algorithm is fed with an input and output dataset of desired solution. The desired output data is otherwise known as labels. These information help train the machine learning model in the precise prediction of output when given a related output after being trained. Some

examples of common algorithms used for supervised learnings are regression analysis - linear regression, etc.-, decision trees, amongst others. [?] The two types of supervised learning techniques used are classification and regression.

- **Unsupervised Learning**

This refers to the machine learning technique in which the algorithm is only fed with input data for the training process. The model learns "unsupervised" by finding out and grouping patterns in inputted datasets. Unlike supervised learning, they require large amount of unlabeled datasets to train -properly find patterns in datasets.

- **Semisupervised Learning**

this is a learning technique that uses both labeled -input and desired output- and unlabeled -only input- datasets for training machine learning models.

The other methods used in training a machine learning model are reinforcement learning and deep learning. Deep learning is a subset of machine learning whose algorithm defines an artificial neural network -ANN- that is designed to emulate the way an human brain learns, and unsupervised or semisupervised learning technique to train [?]. On other hand, reinforcement learning is a reward based learning in the sense that the machine learning is trained by being given a point for reaction to certain events.

### B. Linear Regression Model

Regression is a predictive analysis with a long but glorious history from its successful applications to problems in the statistics and economics domain. Regression is a kind of supervised learning technique for determining the best fit line to describe patterns in data: linear regression uses a straight line to describe these patterns. [?]. The best fit line is the line that reduces the summed squared difference between the value of the line of a certain value  $x$  and its corresponding  $y$  values [?]. The mathematical expression for linear regression is:

$$y = \beta X + \beta_0 \quad (1)$$

Where  $y$  is the dependent variable,  $X$  is the independent variable of the equation,  $\beta$  is a coefficient which represents the slope of the regression line,  $\beta_0$  is a constant value called the **bias**. Equation (1) is the same as the equation of a straight line in linear algebra.

The independent variable  $X$  in (1) is calculated using the mathematical expression below:

$$X = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} \quad (2)$$

Where  $\bar{x}$  and  $\bar{y}$  are the mean values of all the respective  $x$  and  $y$  variables.

Linear regression analysis can be sub-divided into simple and multiple linear regression depending on the number of their independent variables. If the independent variable is one, the regression is said to be a simple linear regression, but if the independent variable is more than one, it is a multiple linear regression. In this paper, we will make analysis using the simple linear regression model.

### C. Model Evaluation

The prediction made by the model using equation (1) can be evaluated using the statistical methods below to find out how close the predicted value is to the actual value.

1)  $R^2$ : this is otherwise known as the coefficient of determination or the coefficient of multiple determination is a measure of how close the best fit line is to the original data using a simple mean. The value of the output of this calculation ranges from 0 to 1 [?]. A value of 1 refers to the fact that all the points fits 100% to the regression line. The lesser the value, the farther the points are to the regression line. The coefficient of determination is defined by the following equation:

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (3)$$

Where  $n$  is the number of observations i.e. total number of variables,  $\hat{y}_i$  is the estimated value of the dependent variable for the  $i^{\text{th}}$  observation computed by the regression equation,  $y_i$  is the observed value of the dependent variable for the  $i^{\text{th}}$  observation and  $\bar{y}$  is the mean of all  $n$  observations of the dependent variable [?]

2) *Mean Squared Error*: Mean Squared Error (MSE) is a performance measure that helps determine how much error is made in the predicted value in relation to the actual value of the output given the same input. The equation of the calculation is denoted by:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \quad (4)$$

where  $x_i$  is a vector of all  $i^{\text{th}}$  observation,  $\hat{f}$  is the prediction function. The lower the value of MSE, the more accurate the prediction is.

## IV. ALGORITHM IMPLEMENTATION

All the algorithm needed for prediction using the linear regression model have been implemented in machine learning libraries like Scikit learn which we will be using in this paper.

### A. Libraries and Tools used

1) *Scikit Learn*: Scikit learn is an open-source python library that provides algorithms that are used in machine learning. This library provide functionalities for solving machine learning jobs like regression, classification, clustering, model selection, pre-processing - like splitting data sets into test and train subsets -amongst others [?]. The scikit-learn API are designed around this main design principles which are consistency - all objects (basic or composite) share a consistent interface composed of a limited set of methods -, inspection - parameters are exposed as public attributes -, composition, sensible defaults - provides understandable default parameters which gives baseline solution for tasks at hand -and nonproliferation of classes - datasets are represented as NumPy arrays or SciPy sparse matrices [?].

2) *Numpy*: NumPy is an open-source Python library that provides routines that allows for fast operations on multidimensional arrays, including mathematical, logical, shape manipulation, sorting, selecting, I/O, discrete Fourier transforms, basic linear algebra, basic statistical operations, random simulation <sup>1</sup>.

3) *Pandas*: Pandas is an open-source data analysis and manipulation python library that allows high-performance easy exploration, cleaning and processing of tabular data structures in Python <sup>2</sup>.

4) *Matplotlib*: Matplotlib is a Python library used along with Numpy for creating static, animated and interactive visualization of data <sup>3</sup>.

5) *Seaborn*: Seaborn is a Python library which builds upon the functionalities of matplotlib and integrates closely with pandas data structures for making statistical graphics to datasets <sup>4</sup>.

## V. PRACTICAL EXAMPLE

In this paper, the california housing datasets from [?] included in the scikit sklearn datasets library is used to predict train and prediction the average prices of housing in california using linear regression.

### A. Preparing Data

The data used in this example is gotten from scikit learn. this returns the following datasets with each attribute meaning....get from this link [https://scikit-learn.org/stable/modules/generated/sklearn.datasets.fetch\\_california\\_housing.html](https://scikit-learn.org/stable/modules/generated/sklearn.datasets.fetch_california_housing.html).

- the libraries were imported
- the dataset was created as a panda data frame to allow easy manipulation I guess. then the target column was added to the dataset columns which will be the y variable
- a description of the data was manage to get a quick statistics on the data being worked on.
- the input variable is the AveRooms ..... it total-room, which the totatl number of rooms in a block,
- We'll randomize the data, just to be sure not to get any pathological ordering effects that might harm the performance of Stochastic Gradient Descent.[google colab]
- use the

### B. Result and discussion

The coefficients are the most important output that we can obtain from our regression model because they allow us to re-create the weighted summation that can predict our outcomes.[regression textbook, page 70]

## VI. CONCLUSION

finish by highlighting the strenth of linear regression as well as the shorcomings.

<sup>1</sup><https://numpy.org/doc/stable/user/whatisnumpy.html>

<sup>2</sup><https://pandas.pydata.org/docs/>

<sup>3</sup><https://matplotlib.org/>

<sup>4</sup><https://seaborn.pydata.org/introduction.html>