# Spark based "linear Regression" & "Sample linear regression"

J Anil kumar - 21BDS024, R Vinay kumar - 21BDS056 , Sidharth Kaushik - 21BDS064

*Abstract*— This report provides an overview of Hadoop, SingleNode, MultiNode, Spark, and Spark Cluster. It also explores Simple Linear Regression, a dataset commonly used in this context, and the Linear Regression algorithm in Spark. Hadoop is a distributed computing system that allows for the storage and processing of large data sets across clusters of computers. SingleNode is a standalone version of Hadoop that runs on a single machine, while MultiNode is a more complex setup that involves multiple machines. Spark is a fast and flexible data processing engine that can be used with Hadoop or as a standalone system. Spark Cluster is a distributed version of Spark that runs on a cluster of computers. Simple Linear Regression is a statistical method used to model the relationship between two variables, while Linear Regression algorithm in Spark is a machine learning technique that uses linear regression to make predictions about a data set. Overall, this report provides a useful introduction to these concepts and techniques for those interested in working with large data sets and distributed computing systems.

## I. INTRODUCTION

Linear regression is a statistical approach used to model the relationship between a dependent variable and one or more independent variables. The goal of linear regression is to find the best-fit line that represents the linear relationship between the variables. This line is represented by the equation $Y = aX + b$, where Y is the dependent variable, X is the independent variable, a is the slope of the line, and b is the y-intercept. The linear regression algorithm calculates the values of a and b that minimize the difference between the predicted values and the actual values.

A linear sample dataset is a set of data that can be used to perform linear regression analysis. It typically consists of two variables, a dependent variable and an independent variable, that are assumed to have a linear relationship. For example, a linear sample dataset could include data on the number of hours studied and the grade received on an exam. The independent variable would be the number of hours studied, and the dependent variable would be the grade received. By analyzing this data with linear regression, we can determine the relationship between the two variables and make predictions about future outcomes.

## II. PROJECT WORK

### A. Benefits of HDFS Cluster for Linear Regression algorithm

The main benefit of using an HDFS cluster for the linear regression algorithm is that it provides a distributed computing framework for the algorithm. This means that the algorithm can be executed in parallel on multiple nodes, allowing for faster processing times and more efficient use of resources. By breaking the data into smaller chunks and distributing them across the HDFS cluster, the algorithm can be executed on each node simultaneously, allowing for faster processing times.

Linear regression is a commonly used statistical technique for predicting a continuous variable based on one or more predictor variables. When performing linear regression on large datasets, the computation can become time-consuming and resource-intensive. This is where HDFS and a distributed computing framework like Apache Hadoop can help.

With HDFS, the dataset can be stored in a distributed manner across multiple nodes in a cluster. This allows the data to be accessed and processed in parallel by multiple nodes, which can greatly speed up the computation time. Additionally, HDFS provides fault tolerance by replicating data across multiple nodes, ensuring that the data is not lost in the event of a node failure.

To perform linear regression using a distributed computing framework like Apache Hadoop, the algorithm can be implemented as a MapReduce job. The MapReduce paradigm is a programming model for processing large data sets in a distributed computing environment. In this approach, the input data is divided into smaller chunks, which are processed in parallel by different nodes in the cluster. The results from the different nodes are then combined to produce the final output.

The MapReduce job for linear regression would involve two phases - the map phase and the reduce phase. In the map phase, each node in the cluster would process a subset of the input data and compute the local linear regression coefficients. In the reduce phase, the results from the different nodes would be combined to produce the final regression coefficients.

Overall, HDFS and a distributed computing framework like Apache Hadoop can help to greatly improve the performance of linear regression on large datasets by enabling parallel processing and fault tolerance.

### B. Linear Regression Algorithm on Sample Linear Regression Data set

Linear regression is a statistical technique used to model the relationship between a dependent variable and one or more independent variables. The input data for a linear regression algorithm typically consists of a set of observations or data points that include values for the dependent variable and one or more independent variables.

The dependent variable, also known as the response variable, is the variable that is being predicted or modeled by the linear regression algorithm. It is denoted by "y" and can be a continuous variable such as a numerical value, or a categorical variable such as a binary value.

The independent variables, also known as the predictor variables, are the variables that are used to predict the value of the dependent variable. They are denoted by "x" and can also be continuous or categorical variables. In simple linear regression, there is only one independent variable, whereas in multiple linear regression, there are multiple independent variables.

The output generated by a linear regression algorithm is a set of coefficients that define the linear relationship between the dependent variable and the independent variables. In simple linear regression, there are two coefficients - the intercept and the slope. The intercept represents the predicted value of the dependent variable when the independent variable is zero, while the slope represents the change in the dependent variable for a unit change in the independent variable.

In multiple linear regression, there are multiple coefficients, each representing the change in the dependent variable for a unit change in the corresponding independent variable, holding all other independent variables constant. The coefficients can be used to make predictions of the dependent variable for new observations or data points.

In addition to the coefficients, the output of a linear regression algorithm typically includes a measure of the goodness of fit of the model, such as the R-squared value. The R-squared value represents the proportion of the variation in the dependent variable that is explained by the independent variables in the model. A higher R-squared value indicates a better fit of the model to the data.

## III. CONCLUSIONS

In conclusion, linear regression is a powerful and widely used algorithm for predicting a continuous target variable based on the relationship with independent variables. It assumes a linear relationship between the features and the target, estimating coefficients that minimize the difference between predicted and actual values. Linear regression provides interpretable results, such as coefficient values that indicate the strength and direction of the relationship. It is important to validate the model's performance using evaluation metrics like mean squared error and R-squared. While linear regression has its limitations, it serves as a fundamental and versatile tool in data analysis and predictive modeling tasks.

## ACKNOWLEDGMENT

References are important to the reader; therefore, each citation must be complete and correct. If at all possible, references should be commonly available publications.

### REFERENCES

[1] https://spark.apache.org/
[2] https://github.com/apache/spark/blob/master/examples/src/main/java/org/apache/spark/examples/ml/JavaKMeansExample.java
[3] https://computingforgeeks.com/how-to-install-apache-spark-on-ubuntu-debian/
[4] https://github.com/apache/spark/blob/master/examples/src/main/java/org/apache/spark/examples/ml/JavaKMeansExample.java