A PROJECT REPORT


ON


# "CALCULATED PREDICTIVE ANALYSIS


# MODEL USING XGBoost AND DECISION TREE"




SUBMITTED BY:

Shivanand Thakur

Master of Computer Application

Vellore Institute of Technology

# Table of Contents

# 1. Introduction

## 1.1 Calculated Predictive Analysis Model an Overview

Stock Market prediction and analysis is the act of trying to determine the future value of a company stock or other financial instrument traded on an exchange. Stock market is the important part of economy of the country and plays a vital role in the growth of the industry and commerce of the country that eventually affects the economy of the country. Both investors and industry are involved in stock market and wants to know whether some stock will rise or fall over certain period of time. The stock market is the primary source for any company to raise funds for business expansions. It is based on the concept of demand and supply. If the demand for a company's stock is higher, then the company share price increases and if the demand for company's stock is low then the company share price decrease.

The **Bombay Stock Exchange** (BSE) is the leading stock exchange of India, located in Mumbai. The BSE was established in 1986 as the first demutualized electronic exchange in the country. BSE became the first company to be popular among the people after NSE as it has lot more companies listed as compared to NSE. BSE has over 5,500 listed companies while NSE had 1,600 listed companies which make BSE more diversified in terms of the companies that can be

The **BSE** index is National Stock Exchange of India's benchmark broad based stock market index for the Indian equity market. It represents the weighted average of 500 Indian company stocks in 12 sectors and is one of the two main stock indices used in India, the other being the NSE, Sensex.

Due to involvement of many numbers of industries and companies, it contains very large sets of data from which it is difficult to extract information and analyse their trend of work manually. The application developed in this project, not only helps in prediction the future movement if the stock in the market, but also automate the data retrieval, trend analysis, predictive analysis and insights generation of a stock, just at the click of a button. Stock market analysis and prediction will reveal the market patterns and predict the time to purchase stock. The successful prediction of a stock's future price could yield significant profit. This is done using large historic market data of 12 months in this project, to represent varying conditions and confirming that the time series patterns have statistically significant predictive power[1].

## 1.2 Objective and Scope of the Project

**Objective:**

> - **To add to the academic understanding of stock market prediction:**
>   - With a greater understanding of how the market moves, investors will be better equipped to prevent another financial crisis.
>   - Evaluate some existing strategies from a rigorous scientific perspective and provide a quantitative evaluation of new strategies.
>
> - **Provide an automated Stock Prediction Tool to Traders to:**
>   - Make Buy/Sell Decisions
>   - Distinguish between conservative and risky stocks

**Scope:**

> - No attempt is made in this project at portfolio management. Portfolio management is largely an extra step done after an investor has made a prediction on which direction any particular stock will move. The investor may choose to allocate funds across a range of stocks in such a way to minimize his or her risk. For instance, the investor may choose not to invest all of their funds into a single company lest that company takes unexpected turn. A more common approach would be for an investor to invest across a broad range of stocks based on some criteria he has decided on before. **This project will focus exclusively on predicting the daily trend (price movement) of individual stocks. The project will make no attempt to deciding how much money to allocate to each prediction. More so, the project will analyse the accuracies of these predictions**[2].
> - A distinction must be made between the trading algorithms studied in this project and high frequency trading (HFT) algorithms. HFT algorithms make little use of intelligent prediction and instead rely on being the fastest algorithm in the market. These algorithms operate in fractions of a second. **The algorithms presented in this report will operate on the order of days and will attempt to be truly predictive of the market.**

## 1.3 Literature Survey

In the last few decades forecasting of stock returns has become an important field of research. In most of the cases the researchers had attempted to establish a linear relationship between the input macroeconomic variables and the stock returns. After the discovery of nonlinearity in the stock market index returns, many literatures have come up in nonlinear statistical modelling of the stock returns, most of them required that the nonlinear model be specified before the estimation is done. But since stock market return is noisy, uncertain, chaotic and nonlinear in nature, Predictive Modelling & Machine Learning has evolved in capturing the structural relationship between a stock's performance and its determinant factors more accurately than many other statistical techniques.

In literature, different sets of input variables are used to predict stock returns. In fact, different input variables are used to predict the same set of stock return data. Some researchers used input data from a single time series where others considered the inclusion of heterogeneous market information and macroeconomic variables. Some researchers even pre-processed these input data sets before feeding it to the Predictive Model for forecasting.

Min and Lee were doing prediction of bankruptcy using machine learning. They evaluated methods based on Support Vector Machine, multiple discriminant analysis, logistic regression analysis, and three-layer fully connected back-propagation neural networks. Their results indicated that support vector machines outperformed other approaches. A Decision Tree is a useful and popular classification technique that inductively learns a model from a given set of data. One reason for its popularity stems from the availability of existing algorithms that can be used to build decision trees, such as CART [3], ID3 [4], and C4.5 [5]. These algorithms all learn decision trees from a supplied set of training data, but do so in slightly different ways. As discussed in the introduction, a classifier is built by analysing training data. That is to say, a classifier is built by analysing a collection of instances where each instance is composed of a set of attribute values and a known class value. These decision tree algorithms build top-down structures that partition instances into separate classes, and it is hoped that these structures generalize well to instances with unknown class values. This would mean that the decision trees have fulfilled their objectives and have indeed discovered some underlying property of the data [4].

Tsai and Wang did research where they tried to predict stock prices by using ensemble learning, composed of decision trees and artificial neural networks. They created dataset from Taiwanese

stock market data, taking into account fundamental indexes, technical indexes, and macroeconomic indexes. The performance of Decision Tree + Artificial Neural Network trained on Taiwan stock exchange data showed Score performance of 77%. Single algorithms showed F-score performance up to 67%.

Josip Arneric, Elza Jurun, Snjezana pivac describes that technical analysis is done to find out the price movements whereas fundamental analysis is done to predict values by looking at the fundamentals of a particular company. They focused on technical analysis and define that trend can be of two types on the basis of either time structure or general direction.

Professor Veroljub says that the way of investing is to sell when prices are at top and to buy when prices are at lower whatever the patterns are. In his articles he has discussed the market efficiency theory, Classical theory, confidence theory and Dow Theory. He also differentiates between the Classical and Confidence theory

Wing-Keung Wong, Meher Manzur and Boon-Kiat Chew (2002) article discuss that the helpful principle of technical analysis is to identify trends and then go with the trend whether it is occurring randomly or due to fundamental factors. He also discussed the techniques of moving averages and relative strength index (RSI) by applying it on Singapore stock exchanges.

There are many tools and software available out there that provide forecasting of stock market entities, share quantity and share value for a given financial organization. Most of them claim to predict the stock market with near to 100% accuracy but the opinions from the users vary. Some of the popular tools and software with their methodologies are mentioned as follows.

➢ **inteliCharts Predictive Stock Market Analytics:**
It is a quantitative modelling tool used for financial time series forecasting. The system is adaptive in its core as it learns the patterns and geometrical relationships defined by historical time series data points, which are unique for each individual stock, index, or another financial instrument.

➢ **Markettrak**
Its stock market forecast system consists of two major parts: an extensive database and a forecast model. The forecast model reads the database and then makes a prediction of where the market is headed. From this prediction, it determines a trading position for the Dow Diamonds or the SP500 Spiders. The database and forecast are updated daily at the close of trading.

- ➢ **Stock-Forecasting.com**

  www.stock- forecasting.com (Centre of Mathematics & Science, Inc., Chicago, United States of America) provides innovative price-prediction technology for active Day Traders, Short- and Long-term Investors. They develop web-based software for stock market forecasting and analysis [2].

The Stock-Forecasting software predicts stock prices, generates trading "Buy-Hold-Sell" signals, computes the most profitable company to invest in and analyses the accuracy of predictions.

The researcher conducted a study of internal representation and a purchasing and selling time prediction system for Tokyo Stock Exchange stocks are covered. Modular neural networks are the foundation of it. For the TOPIX (Tokyo Stock Exchange Prices Indexes) prediction system, we created a number of learning algorithms and prediction techniques. The simulation on stock trading produced outstanding profits, and the prediction system produced precise predictions. Fujitsu and Nikko Securities created the forecasting system[6].

The author of the paper says that intelligent systems may be used to anticipate the stock market. The use of hybridised soft computing approaches for automated trend analysis and stock market forecasting is the topic of this article. We employ a neural network for stock forecasting one day in advance and a neuro-fuzzy system for trend analysis of the anticipated stock values[7]. We took into consideration the well-known Nasdaq-100 index of Nasdaq Stock MarketSM to illustrate the suggested methodology. We looked at 24 months' worth of stock data for the Nasdaq-100 main index as well as six of the Nasdaq100 index businesses. Principal component analysis was used to preprocess the input data before feeding it to an artificial neural network for stock forecasting. A neuro-fuzzy system is further given the expected stock values.

The researcher says that The ways low- and medium-technology (LMT) companies use to create technical innovation are discussed in this research, which also looks at how these strategies affect the innovation performance of the companies. Based on a sample from the Taiwanese Technological Innovation Survey, these analyses a total of 753 LMT companies. According to the descriptive data, 95% of businesses got their technology through technology licencing, whereas 32% of businesses outsourced their R&D. About 20%, 18%, 8%, and 23% of the companies in the sample collaborate with suppliers, clients, rivals, and research organisations to acquire external technological expertise. The results of this study, which use a moderated hierarchical regression analysis, are intriguing. First, the performance of innovation

is not considerably impacted by internal technology licencing. Second, internal R&D spending has a detrimental moderating impact [8].

The Researcher says The genetic quantum algorithm (GQA) is a brand-new evolutionary computing technique that is proposed in this research. The foundation of GQA is the idea and fundamentals of quantum computing, including qubits and superposition of states. Due to the qubit chromosome's probabilistic representation, GQA may express a linear superposition of solutions in place of binary, numeric, or symbolic representation. Quantum gates are used to seek for the optimal answer as genetic operators. The performance of GQA is characterised by quick convergence and strong global search capabilities. Experimental findings on the knapsack problem, a well-known combinatorial optimisation issue, show the efficacy and applicability of GQA. The findings demonstrate that GQA outperforms alternative genetic algorithms that make use of penalty functions, repair techniques and decoder [9].

Given how complicated and volatile the stock market is, experts have long been interested in predicting price movements. stock market intrinsic turbulence throughout the globe makes making predictions difficult. Even if they are useful, forecasting and diffusion modelling cannot be the answer to all of the issues associated with prediction, whether it be long-term or short-term. To ensure low investment risk, market risk, which is closely connected with forecasting mistakes, needs to be reduced. By framing the forecasting problem as a classification problem, which is addressed by a common set of machine learning methods, the authors suggest reducing forecasting error. In this study, we provide a unique approach for reducing the risk of stock market investment through stock return forecasting [10].

The field of software engineering was created in response to the perceived software crisis in the industry. A well, that is. It is a well-known truth that the software business needs to know how much a project will cost to produce and how much time it will take. Compared to resource estimate in other industries, software engineering presents the greatest challenge. The neural network model is one of the resource estimating techniques that are currently accessible. This study suggests comparing several neural network model training techniques and demonstrates which is most appropriate for software engineering applications [11].

# 2. SYSTEM ANALYSIS

## 2.1 Proposed system

### 2.1.1 Defining the Problem

The Stock Market is a complex and dynamical system, & is influenced by many factors that are subject to uncertainty. So, it is a difficult task to forecast stock price movements. Due to technology and globalization of business & financial markets it is important to predict the stock prices more quickly & accurately. Automated User-friendly Trading application can be developed based on financial predictive indicator algorithms & machine learning techniques to predict the performance of stocks in NSE's BSE Index.

**Definition of Problem**

The Stock Market is a complex and dynamical system, & is influenced by many factors that are subject to uncertainty. So, it is a difficult task to forecast stock price movements. Due to technology and globalization of business & financial markets it is important to predict the stock prices more quickly & accurately. Automated User-friendly Trading application can be developed based on financial predictive indicator algorithms & machine learning techniques to predict the performance of stocks in BSE Index.

Investors prefer stock market investments as they have the opportunity of highest return over other schemes. Nifty (benchmark of NSE India) is a well-diversified index consisting of 50 major stocks from 21 sectors of the Indian economy. However, trading through Stock buy/sell prediction computer algorithms is still in its nascent stage in the Indian stock market. The need of an automated user-friendly trading predictor system, which predicts stock price upward/downward movements, is necessary in the Indian stock market, given the explosion of algorithmic trading, being one of the most prominent trends in the global financial industry over recent decade. A Stock Prediction Application will be developed in this project using Nifty data, keeping in mind the following three steps:

1. The system has to have some models generating Stock Market predictions using **Financial Stock Predictor** Functions **(E.g.: Williams %R)** and **Machine Learning** Techniques **(E.g.: Decision Trees)**[12].

2. **Back testing** of the **Prediction Models** is essential to evaluate the trading system's performance on historical market data and thus determine the viability of the system

3. An analytical insight has to be provided, of whether the stock is: **"Bullish"** or **"Bearish".**

## 2.1.2 Developing Solution Strategies

Developing a stock market prediction system using machine learning involves several steps and strategies. Here's a general outline of the process:

Data Collection: Gather historical stock market data, including price, volume, and other relevant indicators. You can obtain this data from various financial data providers or APIs.

Data Preprocessing: Clean and preprocess the collected data to handle missing values, outliers, and inconsistencies. Perform tasks like data normalization, feature scaling, and handling categorical variables.

Feature Selection: Identify the most relevant features that have a significant impact on stock price movements. This can be done through techniques like correlation analysis, feature importance, or domain expertise.

Feature Engineering: Create additional features that can capture meaningful patterns or trends in the data. For example, you can calculate moving averages, technical indicators (e.g., RSI, MACD), or derive features from news sentiment analysis.

Model Selection: Choose an appropriate machine learning algorithm for your prediction task. Commonly used algorithms for stock market prediction include linear regression, decision trees, random forests, support vector machines (SVM), or more advanced techniques like deep learning models (e.g., recurrent neural networks or long short-term memory networks).

Model Training: Split the pre-processed data into training and testing sets. Train the selected model using the training data while adjusting the model's hyperparameters. Cross-validation techniques like k-fold cross-validation can be used to optimize the model's performance and prevent overfitting.

Model Evaluation: Assess the performance of your trained model using appropriate evaluation metrics such as mean squared error (MSE), root mean squared error (RMSE), mean absolute error (MAE), or accuracy metrics depending on the nature of the problem (regression or classification).

Model Refinement: Fine-tune your model by iterating on the previous steps. Experiment with different algorithms, hyperparameters, or feature combinations to improve the predictive performance. Consider techniques like ensemble learning (combining multiple models) to enhance accuracy.

Back testing and Validation: Apply your trained model to a separate validation dataset or historical data to evaluate its performance in real-world scenarios. Back testing can help simulate how the model would have performed in the past.

Deployment and Monitoring: Once you have a well-performing model, deploy it to a production environment. Continuously monitor its performance and update the model periodically to adapt to changing market conditions and incorporate new data.

It's important to note that stock market prediction is a highly complex and volatile domain, and accurate predictions are challenging to achieve consistently. Financial markets are influenced by various factors, including economic indicators, political events, investor sentiment, and more. Therefore, it's crucial to set realistic expectations and consider the limitations of any prediction system.

## 2.1.4 Data Flow Diagram

A data flow diagram shows the way information flows through a process or system. It includes data inputs and outputs, data stores, and the various sub processes the data moves through. DFDs are built using standardized symbols and notation to describe various entities and their relationships.

**Data flow diagram levels:**

Data flow diagrams are also categorized by level. Starting with the most basic, level 0, DFDs get increasingly complex as the level increases.

**Context Diagrams/Level 0 DFDs** are the most basic data flow diagrams. They provide a broad view that is easily digestible but offers little detail.

**Context Diagram:**

**Level -0 DFD:**



Fig 2: Level 0- Basic Data Flow of the System.

**Level 1 DFDs** go into more detail than a Level 0 DFD. In a level 1 data flow diagram, the single process node from the context diagram is broken down into sub processes. As these processes are added, the diagram will need additional data flows and data stores to link them together.

**Level – 1 DFD:**



Fig 3: Level 1- Data Retrieval & Transformation

## 2.2 Dataset Descriptions

Dataset Description: Calculated Predictive Analysis Model

Dataset Name: BSE Dataset

Dataset Overview: Calculated Predictive Analysis Model Dataset is a collection of data gathered for the purpose of analysing and understanding various aspects related to fine the Stock price of the Companies listed in the BSE Index. It contains a diverse range of data points, including the Date on which the Stock Price movement has taken place, High_price, Low_price, opening_price, Closing_price, No. of share, No. of trades, Total Turn Over, etc,

Data Sources: The dataset is sourced from BSE India (Bombay Stock Exchange) which is the oldest and one of the major stock exchanges in India. It is located in Mumbai and serves as a platform for trading equities, derivatives, debt instruments, mutual funds, and other financial products. The sources include the data of the listed companies on the basis of their daily trades, Monthly trades and yearly trades. It also has the information about the company's Equity, Mutual Funds, Debt and others.

Data Format and Structure: The dataset is provided in CSV file format. It is organized into Rows and Columns. Each record within the dataset represents a distinct information about the date on which the trade has taken place, the opening and the closing price of the trade. It also provides the information about the High_price and Low_price that the stock showed on the specific date. The respective data set of the company has the information about the total number of shares used in the trade that day and even the total number of shares involved in the process.

Variables and Features: The dataset comprises several variables or features that capture relevant information for the project. Some of the key variables includes total turnover of the company, Deliverable Quantity, % Deli. Qty to Traded Qty, Spread High-Low, Spread Close-Open

Data Preprocessing: The dataset has undergone certain preprocessing steps to ensure its quality and usability. This includes removal of all the null and the black data. The dataset is processed in such a way that on the days when the market was closed the data is not left blank at those days (data has be removed of those dates), the dates has been mentioned in a way that the software can read it easily to find the quarterly data of the companies. The data is present in the form of rows and columns therefore it is easy to find the total number of data available in our dataset.

Data Size: The dataset consists of 100gb of records.

Data Usage and License: The Project Calculated Predictive Analysis Model Dataset is intended for research and analysis purposes related to the specific project. It is important to note that the dataset is subject to BSE stock price.

Related Metadata: Alongside the dataset, supplementary metadata is provided to aid in understanding the data. This may include data dictionaries, codebooks, or additional documentation that explains the variables, data collection process, or any transformations applied.

Data Availability: The dataset is available BSE Stock Price. Instructions for accessing or obtaining the dataset can be found at:[21]

## 2.3 Coding Module (Modular Descriptions)

### 1. Importing Libraries



Fig 8: Importing Libraries

### 2. Importing Dataset

The dataset we will use here to perform the analysis and build a predictive model is Reliance Stock Price data. We will use OHLC ('Open', 'High', 'Low', 'Close') data from 1st January 2013 to 31st December 2022 which is for 9 years for the Reliance Power.



Fig 9: Reliance Stock Data

## 3. Counting the total number of Rows and Columns



Fig 10: Total Numbers of Rows and columns

## 4. Describing the Dataset

Describing the data on the basis of count value, mean value, std value, min value, 25% of the data, 50% of the data, 75% of the data and max value.



Fig 11: Describing the Data

## 5. Information of the Data type in the Dataset

Defining the data type of the available dataset.



Fig 12: Data Type Discription

15

## 6. Exploratory Data Analysis

EDA is an approach to analysing the data using visual techniques. It is used to discover trends, and patterns, or to check assumptions with the help of statistical summaries and graphical representations.

While performing the EDA of the Reliance Stock Price data we will analyse how prices of the stock have moved over the period of time and how the end of the quarters affects the prices of the stock.



Fig 13: Trend Analysis

The prices of Reliance Power stocks are showing a mix trend of upward trend and downward trend as depicted by the plot of the closing price of the stocks.

## 7. Head Description for similar data



Fig 14: Determining identical data

If we observe carefully, we can see that there is no similar data that we have generated after our calculation.

## 8. Checking for the availability of null values



Fig 15: Confirming Null values

This implies that there are no null values in the data set provided.

## 9. Creating SubPlots



Fig 16: Subplot for the Open, High, Low, Close, WAP, No.shares

In the distribution plot of OHLC data, we can see two peaks which means the data has varied significantly in two regions. And the No. of shares data is left-skewed.

## 10. Creating BoxPlot



Fig 17: BoxPlot for the Open, High, Low, Close, WAP, No.shares

From the above boxplots, we can conclude that only volume data contains outliers in it but the data in the rest of the columns are free from any outlier.

## 11. Feature Engineering

Feature Engineering helps to derive some valuable features from the existing ones. These extra features sometimes help in increasing the performance of the model significantly and certainly help to gain deeper insights into the data.



Fig 18: Deriving the data of Month, Day, Year

Now we have three more columns namely 'day', 'month' and 'year' all these three have been derived from the 'Date' column which was initially provided in the data.

## 12. Determining whether it is a Quarter end or not using Boolean values



Fig 19: Determining Quarter End

A quarter is defined as a group of three months. Every company prepares its quarterly results and publishes them publicly so, that people can analyse the company's performance. These quarterly results affect the stock prices heavily which is why we have added this feature because this can be a helpful feature for the learning model.

## 13. Creating Bar graph



Fig 20: Bar Chart for the Reliance Data

From the above bar graph, we can conclude that the stock prices increased by 50% from year 2021 to that in 2022.

## 14. Calculating Quarter End



Fig 21: Price Differences of Quarter End

Here are some of the important observations of the above-grouped data:

Prices are higher in the months which are quarter end as compared to that of the non-quarter end months.

The WAP of trades is lower in the months which are quarter end.

## 15. Creating Pie-Chart

Above we have added some more columns which will help in the training of our model. We have added the target feature which is a signal whether to buy or not we will train our model to predict this only. But before proceeding let's check whether the target is balanced or not using a pie chart.



Fig 22: Creating Pie Chart for the Reliance Data

When we add features to our dataset, we have to ensure that there are no highly correlated features as they do not help in the learning process of the algorithm.

## 16. Creating Heatmap

```
In [23]:  plt.figure(figsize=(10, 10))

          # As our concern is with the highly
          # correlated features only so, we will visualize
          # our heatmap as per that criteria only.
          sb.heatmap(df.corr() > 0.9, annot=True, cbar=False)
          plt.show()
```



Fig 23: Heat map determining various values

From the above heatmap, we can say that there is a high correlation between OHLC that is pretty obvious, and the added features are not highly correlated with each other or previously provided features which means that we are good to go and build our model.

## 17. Data Splitting and Normalization



Fig 24: Data Splitting and Normalization

After selecting the features to train the model on we should normalize the data because normalized data leads to stable and fast training of the model. After that whole data has been split into two parts with a 90/10 ratio so, that we can evaluate the performance of our model on unseen data.

## 18. Model Development and Evaluation

Now is the time to train some state-of-the-art machine learning models (Logistic Regression, Support Vector Machine, XGBClassifier), and then based on their performance on the training and validation data we will choose which ML model is serving the purpose at hand better.

For the evaluation metric, we will use the ROC-AUC curve but why this is because instead of predicting the hard probability that is 0 or 1, we would like it to predict soft probabilities that are continuous values between 0 to 1. And with soft probabilities, the ROC-AUC curve is generally used to measure the accuracy of the predictions.



Fig 25: Training Accuracy and Validation Accuracy Value after pridiction

Among the three models, we have trained XGBClassifier has the highest performance but it is pruned to overfitting as the difference between the training and the validation accuracy is too high. But in the case of the Logistic Regression, this is not the case.
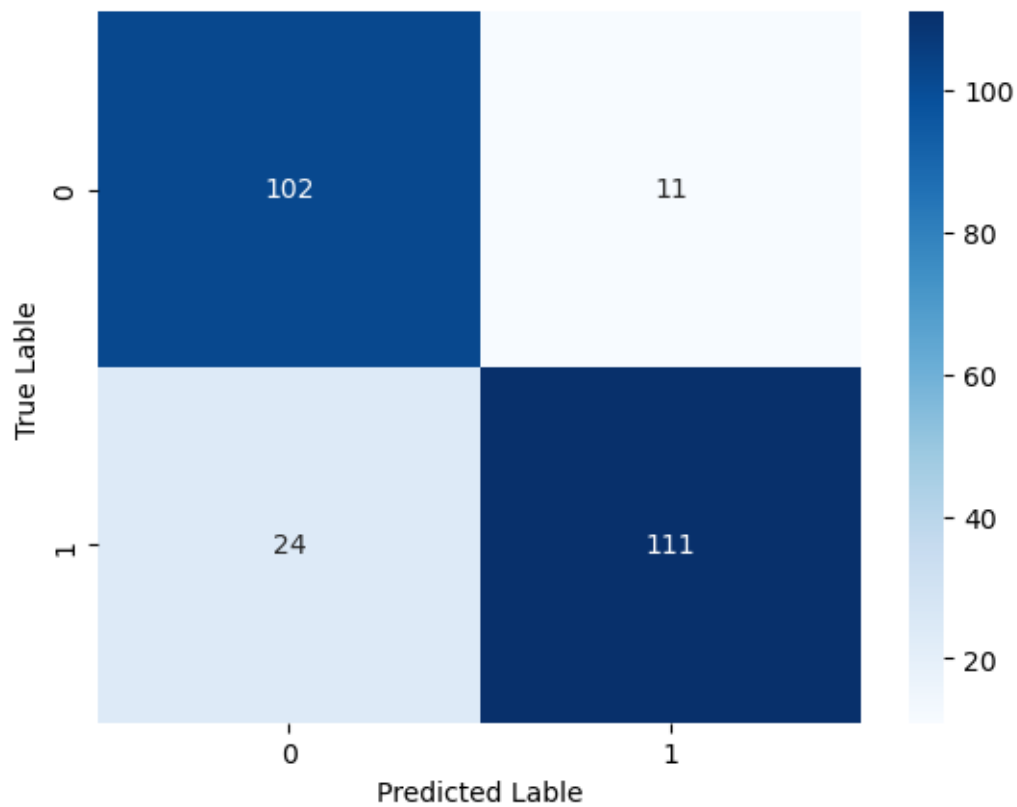
## 19. Confusion Matrix



Fig 26: Creating the Confusion matrix for the Reliance Stocks
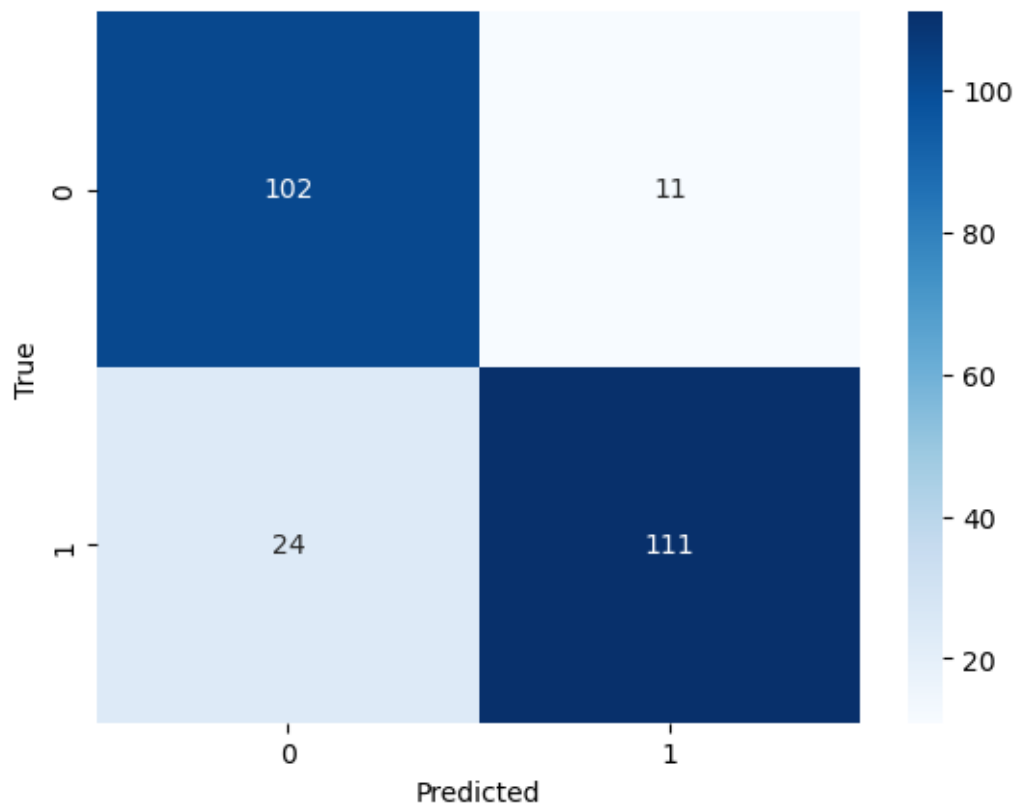
## 20. Calculated Confusion Matrix



Fig 27: Calculated Confusion matrix for the Reliance Stocks

The confusion matrix is a table that shows the performance of a classification model by comparing the predicted labels with the true labels. It provides insights into the number of true positives, true negatives, false positives, and false negatives.

## 21. Calculating the Root Mean Square and Accuracy of the Model



Fig 28: Value of Root Mean Square

## Conclusion

The provided code snippet demonstrates the use of a linear regression model to predict the closing prices of a stock based on given features. Here's a conclusion based on the code:

Data Preparation: The code reads the data from a CSV file named 'Reliance Power.csv' using pandas and splits it into features (X) and the target variable (y). The features include 'Open Price', 'High Price', 'Low Price', 'WAP', and 'No.of Shares'.

Train-Test Split: The data is further split into training and testing sets using the train_test_split function from scikit-learn. The testing set size is set to 20% of the data, and a random state of 2022 is used for reproducibility.

Feature Scaling: The features in the training and testing sets are scaled using the StandardScaler from scikit-learn. The scaling is applied to normalize the feature values.

Model Training: A linear regression model is instantiated using the LinearRegression class from scikit-learn. The model is trained on the scaled training data using the fit method.

Prediction and Evaluation: Predictions are made on the scaled testing set using the trained model, and the mean squared error (MSE), mean absolute error (MAE), and Root Mean-squared (r2) scores are calculated to evaluate the model's performance. The evaluation metrics are printed to the console.

Prediction on New Data: Additionally, the code demonstrates making predictions on new, unseen data. Two rows of data are read from the 'Reliance Power.csv' file, containing the same features used for training the model. The new data is then scaled using the same scaler, and the model predicts the corresponding closing prices. The predicted prices are printed to the console.

In conclusion, the code trains a linear regression model on historical stock data and evaluates its performance using various metrics. It also shows how to use the trained model to predict the closing prices of new, unseen data.
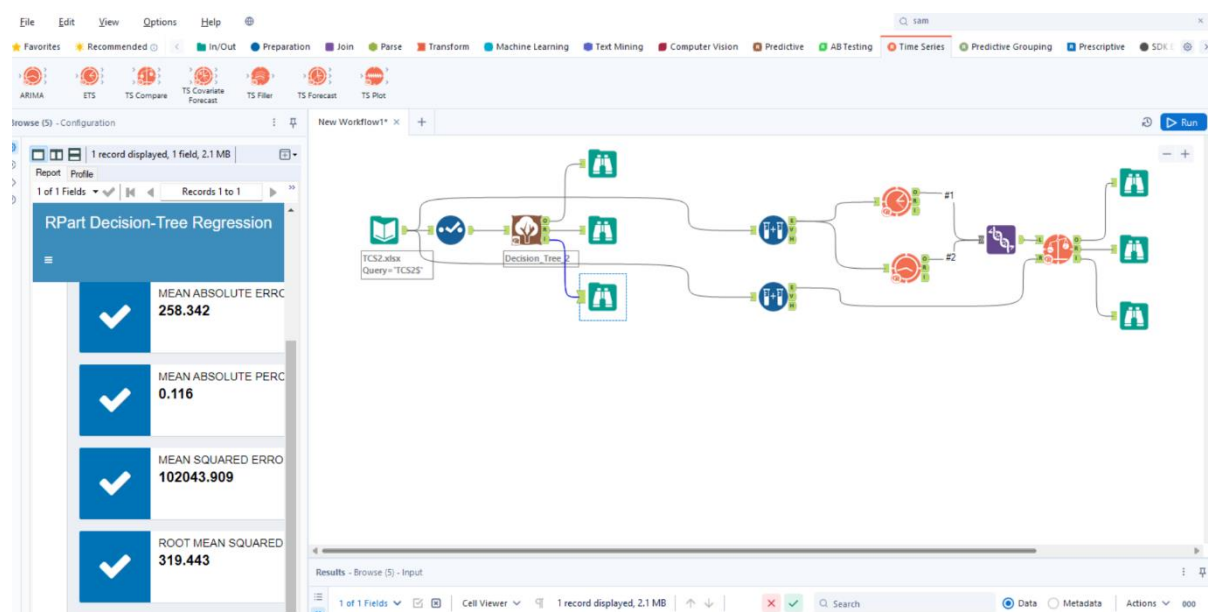
## 22. Comparison With Existing Models



Fig 29: System Design for Model Comparison

Using the Alteryx Software, we Designed a system which is being used to compare our Decision Tree based model with the existing model Arima and ETS.

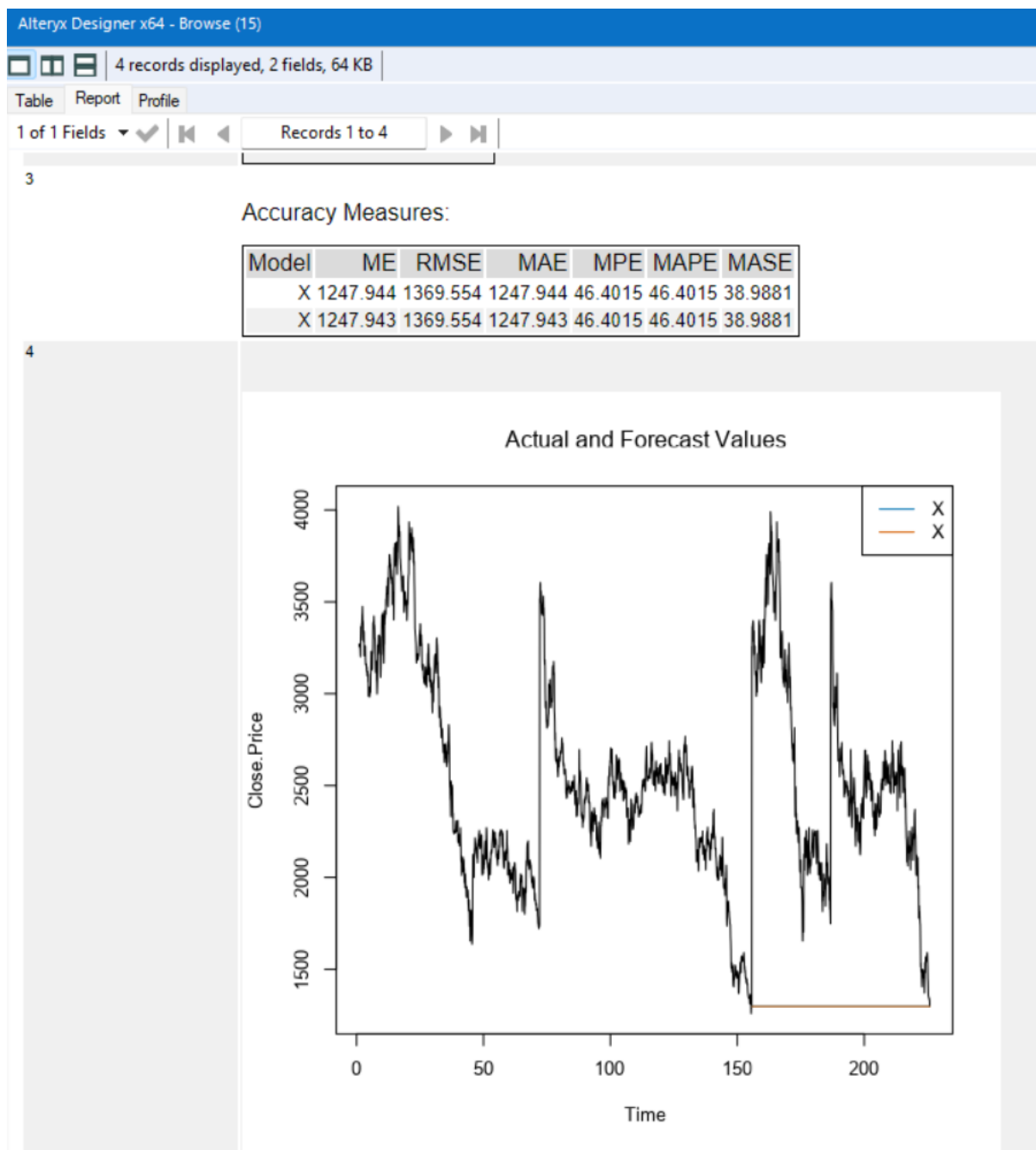## 23. RMS, MAE and ME Values of ARIMA & ETS Model



Fig 30: RMS, MEA & ME value of Arima and ETS Model

AS we can observe in the above output the RMS (Root Mean Square), MAE and ME value for the ARIMA and ETS model is 1369 for the given dataset on which the system is constructed.

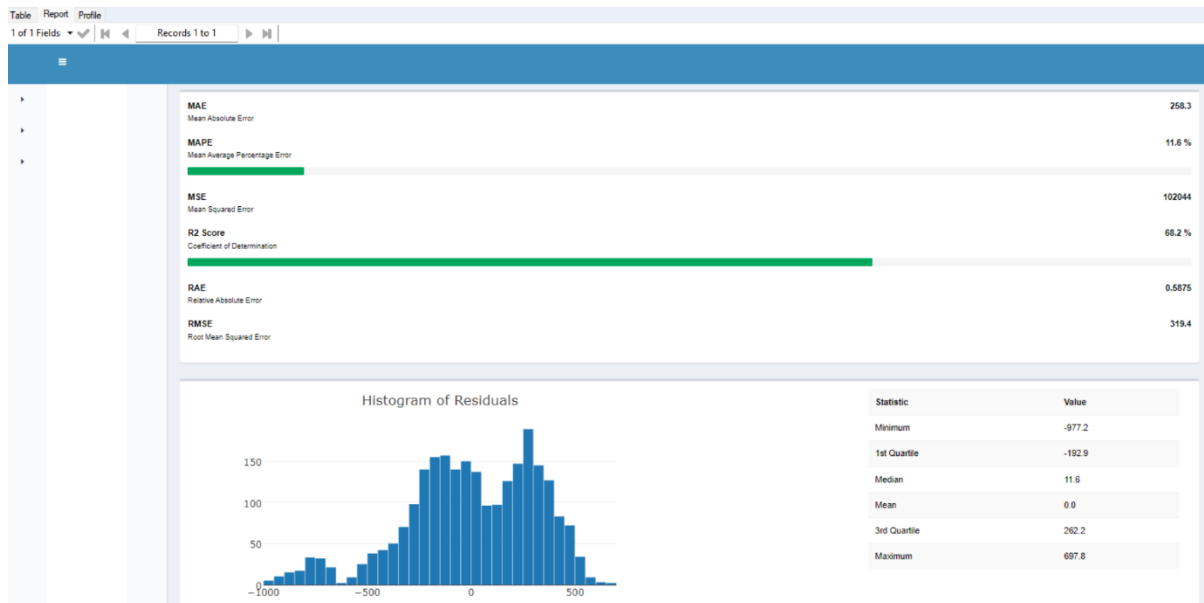## 24. RMS, MAE & ME Values of Decision Tree Model



Fig 31: RMS, MAE & ME value of Decision Tree Model

AS we can observe in the above output the RMS (Root Mean Square), MAE and ME value for the Decision Tree model is just 319 for the given dataset on which the system is constructed which very less when compared with the existing models ARIMA and E

# References

[1] W. K. Wong, M. Manzur, and B. K. Chew, "How rewarding is technical analysis? Evidence from Singapore stock market," *Applied Financial Economics*, vol. 13, no. 7, pp. 543–551, Jul. 2003, doi: 10.1080/0960310022000020906.

[2] Milind Paradkar, "Machine Learning Application in Forex Markets," *quantinsti.com*, Mar. 28, 2016.

[3] M. Segal, "Tree Depth in a Forest NUS / IMS Workshop on Classification and Regression Trees."

[4] J. R. Quinlan, "Induction of Decision Trees," 1986.

[5] by J. Ross Quinlan, M. Kaufmann Publishers, and S. L. Salzberg, "Programs for Machine Learning," 1994.

[6] T. Kimoto, K. Asakawa, M. Yoda, and M. Takeoka, "Stock Market Prediction System with Modular Neural Networks."

[7] A. Abraham, B. Nath, and M. P. K, "Hybrid Intelligent Systems for Stock Market Analysis."

[8] K. H. Tsai and J. C. Wang, "External technology sourcing and innovation performance in LMT sectors: An analysis based on the Taiwanese Technological Innovation Survey," *Res Policy*, vol. 38, no. 3, pp. 518–526, Apr. 2009, doi: 10.1016/j.respol.2008.10.007.

[9] K.-H. Han, "Genetic Quantum Algorithm and its Application to Combinatorial Optimization Problem."

[10] L. Khaidem, S. Saha, and S. R. Dey, "Predicting the direction of stock market prices using random forest," Apr. 2016, [Online]. Available: http://arxiv.org/abs/1605.00003

[11] K. K. Aggarwal, Y. Singh, P. Chandra, and M. Puri, "Evaluation of various training algorithms in a neural network model for software engineering applications," *ACM SIGSOFT Software Engineering Notes*, vol. 30, no. 4, pp. 1–4, Jul. 2005, doi: 10.1145/1082983.1083003.

[12] Chainika Thakar, "Machine Learning for Trading," *quantinsti.com*, Aug. 16, 2021.

[13] Larry Page and Sergey Brin, "Google Finanace," *google.com*, Sep. 04, 1998.

[14] "trading algorithms work," *quora.com*.

[15] QuantInsti Blog, "Algo Trading & Quant Finance," *quantinsti.com*.

[16]    andredumas, "Data Plots,Indicators,Interaction & Other Examples," *github.com*, Sep. 28, 2016.

[17]    Chainika Thakar, "Regression," *quantinsti.com*, Apr. 18, 2023.

[18]    Sushant Ratnaparkhi & Milind Paradkar, "Decision Trees," *quantinsti.com*, Oct. 17, 2017.

[19]    Eric Hammer, "Machine Learning for Quantitative Finance," *quantinsti.com*, Apr. 28, 2017.

[20]    "Financial companies use machine learning," *www.quora.com*.

[21]    Premchand Roychand, "BSE Index," *bseindia.com*, Jul. 09, 1875.