

Prediction of Stock Prices using Random Forest and Support Vector Machines

S. Arun Kumar, Abhishek Jha, Shashank Shekhar, Ashutosh Kumar Singh

Abstract: Prediction of stock markets is the act of attempting to determine the future value of an inventory of a business or other financial instrument traded on an economic exchange. Effectively foreseeing the future cost of a stock will amplify the benefits of the financial specialist. This article suggests a model of machine learning to forecast the price of the stock market. During the way toward considering various techniques and factors that should be considered, we found that strategy, for example, random forest, support vector machines were not completely used in past structures. In this article, we will present and audit an increasingly suitable strategy for anticipating more prominent exactness stock oscillations. The primary thing we thought about was the securities exchange estimating informational index from yahoo stocks. We will audit the utilization of random forest after pre-handling the data, help the vector machine on the informational index and the outcomes it produces. The powerful stock gauge will be a superb resource for financial exchange associations and will give genuine options in contrast to the difficulties confronting the stock speculator.

Keywords : Multiple instance learning, Support Vector Machine, Random forest, data set, stock market.

I. INTRODUCTION

Basically, the stock market is an aggregation of different inventory buyers and vendors. A stock generally reflects a person or a group of people's ownership claims on company. The attempt to determine the stock market's future value is known as a forecast will be robust, precise and effective. The framework should work as per the genuine situations and be great appropriate for true setups. It is also expected that the system will consider all the variables that could affect the value and performance of the stock. There are different techniques and ways to implement the forecast scheme such as fundamental analysis, technical analysis, machine learning, market mimicry, and aspect structuring of the time series. The forecast has shifted into the technological domain with the advance of the digital era. Using Artificial Neural Networks, Recurrent Neural Networks, which is basically AI implementation, is the most noticeable and promising method. Machine learning includes artificial intelligence that empowers the system to learn from previous experiences and enhance them without being programmed

Revised Manuscript Received on November 15, 2019

Mr. Arun Kumar, Asst. Professor, Computer Science and Engineering, SRM Institute of Science and Technology, Ramapuram, Chennai, India

Abhishek Jha, B.Tech Student, Computer Science and Engineering, SRM Institute of Science and Technology, Ramapuram, Chennai, India

Shashank Shekhar, B.Tech Student, Computer Science and Engineering, SRM Institute of Science and Technology, Ramapuram, Chennai, India

Ashutosh Kumar Singh, B.Tech Student, Computer Science and Engineering, SRM Institute of Science and Technology, Ramapuram, Chennai, India

over and over. Conventional forecast methods in AI use calculations, for example, Backward propagation, otherwise called back propagation mistakes. Many scientists have been using more of the teaching methods of the ensemble recently. To predict future highs, it would use low cost and time lags while another network would use lagged highs to predict future highs. Using these projections, inventory prices were performed. The stock price forecast appears to be a random process for short-term windows. The stock price fluctuation generally creates a linear curve over a lengthy period of time. People tend to purchase stocks whose prices in the near future are anticipated to increase. The stock market uncertainty prevents individuals from investing in stocks. Thus, the stock market that can be used in a real-life situation needs to be predicted correctly. The procedures used to foresee the stock market incorporate anticipating a period arrangement alongside specialized examination, demonstrating AI and anticipating the stock market variable. The inventory market prediction model data sets include information such as the closing price, opening price, information and other factors required to predict the item variable which is the price in a specified day. The prior model used traditional forecast techniques such as multivariate analysis with a time series model forecast. Prediction on the stock market surpasses when handles as a regression issue, but when handles as a classification, it works well. The objective is to develop a model that uses machine learning approaches to gain market data and measure future trends in the growth of stock value for both classification and regression. The Support Vector Machine (SVM) can be used for this purpose. These have been noted to be utilized more in portrayal based issues, for example our own. The vector machine method, that we have plotted in which each person statistics phase signify as a point in n-dimensional plane with the estimation of the aspect being the estimation of a specific association and as a result the gathering is executed with the aid of discovering the unambiguously seping hyper plane. For the same, prescient procedures, for example, random forest technique are utilized. The random forest calculation pursues a grouping and relapse learning approach for the ensemble. The random forest requires the normal of the informational index's diverse sub tests, which improves prescient exactness and diminishes informational collection over fitting.

II. PROBLEM DEFINITION

Basically, stock market prediction is described as attempting to determine stock value and providing individuals with a solid concept of knowing and predicting the market and stock prices.

Prediction of Stock Prices using Random Forest and Support Vector Machines

It is usually provided using the data set using the quarterly financial ratio. In this manner, depending on a solitary informational collection may not be sufficient for the conjecture and may offer a wrong result. Therefore, in order to predict the market and inventory trends, we are considering the study of machine learning with different data sets inclusion.

If there is no suggestion of a better stock market prediction algorithm, the problem of stock price estimation will remain a problem. It's quite hard to predict how the stock market will perform. The stock market movement is generally determined by thousands of investors' feelings. Stock market prediction requires an capacity to predict the impact on investors of latest occurrences. These occurrences can be political occurrences such as a political leader declaration, a piece of scam news, etc.

It can also be an global occurrence such as sharp currency and commodity movements, etc. All these occurrences have an impact on corporate income, which in turn impacts investor sentiment.

It is beyond the reach of nearly all investors to predict these hyper parameters properly and consistently. All of these variables make prediction of inventory prices very hard.

III. LITERATURE SURVEY

1. The Stock Market Survey using SVM

Recent studies provide well-founded proof that in sample predictability testing most models of predictive regression are inefficient. The reason for this inefficiency was unstable parameters and uncertainty of the model. The surveys also found the traditional approaches promising to fix this issue. The vector support machine frequently referred to as SVM offers the solution's kernel, decision function, and sparsity. There are many algorithms on the market, but greater performance and precision is provided by SVM. The SVM and stock market correlation analysis shows a powerful interconnection between inventory prices and market index.

2. Impact of Technical Analysis on Stock Price Prediction

An growing trend is the use of machine learning and artificial intelligence methods to predict inventory prices. More and more scientists are investing their time each day in finding methods to arrive at methods that can further enhance the inventory prediction model's precision. Even if the same data set is implemented, the output differs for each method. The quoted article used the stock price forecast using the random forest algorithm to estimate the stock price using the past quarter's financial ratios. This is just one way to look at the issue by using a predictive model to approach it, using the random forest to predict the future stock price from historical information. However, there are always other variables influencing stock price, such as investor feelings, corporate public opinion, news from different outlets, and even events that cause the entire stock market to fluctuate. The precision of the stock price forecast model can be improved by using the financial ratio together with a model that can efficiently evaluate feelings.

3. Corporate Communication Network and Stock Price Movements

This article attempts to show that patterns of communication can have a important impact on the performance of an organization. This article suggested a method to reveal a company's performance.

The method implemented in the journal is used to determine the interactions between important employee email exchange frequencies and business results reflected in inventory values. This article suggested using a data mining algorithm on a publicly accessible Enron Corp data set to detect association and non-association interactions. The Enron Corporation was a Houston, Texas-based power, goods and utilities corporation whose inventory data set is accessible for government use.

4. Machine Learning Approach In Stock Market Prediction

The extensive majority of stockbrokers, while making the prediction, used the particular, indispensable or time collection analysis. Overall, these methods could not be fully trusted, so the need to give economic exchange prediction a powerful approach arose. The methodology chose to be implemented as machine learning and AI together with supervised classifier in order to find the best precise outcome. Results on the binary classification using SVM classifier with an alternative set of a function list were tested. Most of the business-care strategy to machine learning problems benefited from factual methods that excluded AI, despite the reality that there was an optimal method for particular problems. Predicting inventory price used parse records to calculate the prediction, send it to the user, and perform tasks autonomously using the concept of automation, such as buying and selling stocks. Used was the Naïve Bayes algorithm.

5. Stock market prediction using historical analysis.

The stock market forecast technique is loaded up with vulnerability and it very well may be affected by various elements. Hence, in business and fund, the stock market assumes a huge job. The specialized and key investigation is brought out through the sentimental assessment method. Due to its enhanced use, social media information has a strong effect and can be useful in anticipating the stock market trend. Using computer studying algorithms on historic inventory rate information, technical analysis is performed. The connection between different data points is regarded and on these information points a forecast is made. The model has been able to predict future stock values.

IV. DISADVANTAGES OF THE EXISTING SYSTEM

- Most MIL methodologies work in cluster or disconnected mode and thusly accept access to the entire preparing set.
- This limits their relevance in successive and dynamic situations where the data comes.
- Compared to contemporary ML algorithms, MIL is comparatively slow.
- It may be regarded as too strict the normal hypothesis, which may be its disadvantage. The MIL isn't too precise.

V. PROPOSED SYSTEM

We focus on foreseeing stock qualities in this proposed plan utilizing machine learning calculations, for example, Random Forest and Vector Machines Support. We recommended the "Stock Market value expectation" plot we utilized the random forest calculation to foresee the securities exchange cost. We had the option to prepare the machine from the distinctive data focuses from the past in this proposed plan to make a future conjecture. We used information from the stocks of the past year to train the model.

To fix the issue, we mainly used two machine-learning libraries. The first was numpy, used to clean and manipulate the information and put it into a ready-to-analyze form. The other was a sci-kit used to analyze and predict real things.

The informational index that we utilized was stock price from earlier years gathered from the online open database, 80 percent of information was utilized to prepare the machine, and the staying 20 percent were utilized to test the information.

The supervised learning model's basic approach is to know from the practice set the patterns and interactions in the information and then replicate them for the test information. For data processing, we used the python pandas library that combined various data sets into a data frame.

The updated data frame enabled us to prepare the information for extraction of the function. The date and the closing price for a specific day were the data frame characteristics. We used all these features to train the machine on the model of random forests and predicted the variable object, which is the price of a given day. We also quantified the precision using the test set projections and the real values. The suggested scheme affects various study fields, including pre-processing information, random forest, etc.

VI. METHODOLOGIES

1. Random Forest Algorithm

A random forest calculation is used to calculate stock market expectations. Because it is called one of the easiest machine learning calculations to comprehend and adapt, it provides wonderful prediction accuracy. This is generally used in grouping commitments. Because of the high volatility in the stock market, predicting work is quite hard.

We use random forest classifier in stock market prediction that has the same hyper parameters as a decision tree. The decision tool has a comparable model to a tree. It requires the choice based on possible implications, including factors such as the outcome of events, resource costs, and usefulness. The random forest calculation speaks to a calculation wherein it arbitrarily chooses separate perceptions and highlights to manufacture different choice trees and afterward assumes control over the total results from different choice trees. The data is isolated into segments dependent on name or attribute issues. The informational collection we utilized was gathered from the earlier year's securities exchanges from the online open database, 80% of the information was utilized to prepare the machine, and the staying 20% was utilized to test the information. The supervised learning model's basic approach is to know from the practice set the patterns and interactions in the information and then replicate them for the test information.

2. Support Vector Algorithm

The primary job of the supporting machine algorithm is to define an N-dimensional space that categorizes the information points differently. N stands for a number of characteristics here.

There may be various possible hyper planes that can be selected between two classes of information points. The objective of this algorithm is to find a maximum margin plane. Maximizing margin relates to the distance between the two classes' information points. The benefit of maximizing the margin is that it provides some reinforcement to make it simpler for future data points to be classified. Hyper planes are called decision boundaries that assist classify information points. They are ascribed to distinct classes based on the position of the information points relative to the hyper plane. The hyper plane model is based on the amount of characteristics, if there are two characteristics, the hyper plane is a row, if there are three characteristics, then the hyper plane is two-dimensional.

VII. SYSTEM ARCHITECTURE

First step : Extraction of feature is the method of reducing the original set of enormous raw information to more manageable processing groups.

Large set of information involve the processing of a lot of computing resources.

We took information set from kaggle, an online community, to provide free data sets to analyze information.

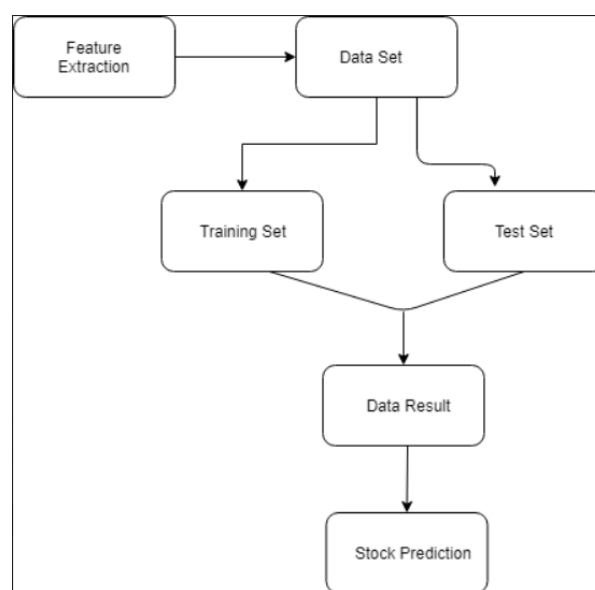
Second step : Classification determines which category belongs to the information set.

Training information for training an algorithm is used.

Test information are used to evaluate model output and enhance precision.

We generally divided the information between testing and training phases around 20 percent -80 percent, and the training data has a greater raw information ratio.

Third step : To predict the stock market price, we use qualified information set results.



VIII. MODULE IDENTIFICATION

Data Collection :

The method of collecting and measuring data from numerous distinct sources is data collection. In order to use the information we collect to create practical alternatives for artificial intelligence and machine learning, it must be gathered and stored in a manner that makes sense for the company at hand.

Our information comprises primarily of preceding year inventory prices from yahoo finance, apple stocks and Kaggle.

Pre-processing :

Pre-processing identifies with the changes that are applied to our data before the calculation is fed. Pre-processing data is a part of data mining that incorporates changing crude data into a progressively steady position. Information Pre-processing is a strategy for changing over crude data into a smooth arrangement of data. In other words, it is collected in raw format whenever the information is obtained from distinct sources, which is not possible for assessment. To use the information we collect to create practical alternatives for artificial intelligence and machine learning, it needs to be collected and stored in a way that makes sense to the company in hand.

Our information comprises primarily of preceding year inventory prices from yahoo finance, apple stocks and Kaggle.

Training the Machine :

The preparation sets are utilized to fit the models. The preparation of the model incorporates cross-approval where we get the inexact presentation of the model utilizing training information. The way toward preparing a ML model includes giving a ML calculation (that is, the learning calculation) with preparing information to gain from.

Data Scoring :

Scoring is additionally called forecast, and is the way toward producing esteems dependent on a prepared AI model, given some new input information. The method used to process the dataset is random forest calculation. Random forest calculation is utilized for arrangement and relapse. The qualities or scores that are made can speak to forecasts of future qualities, yet they may likewise speak to a possible class or result.

IX. EXPERIMENTAL RESULTS

The csv document contains the information which we have used to prepare the machine. It includes traits and sections that shows rise and fall in the financial exchange value forecast. Among these qualities, some of them demonstrates a high worth which demonstrates the most elevated an incentive in the stock that was in the earlier year and some of them demonstrates a low worth which shows the least estimation of past year. OPEN speaks to the beginning day of the exchange and close speaks to the day just before the exchange was closed. VOLUME shows the quantity of offers exchanged on that specific day. NAME speaks to the name of the organization from which the informational index was

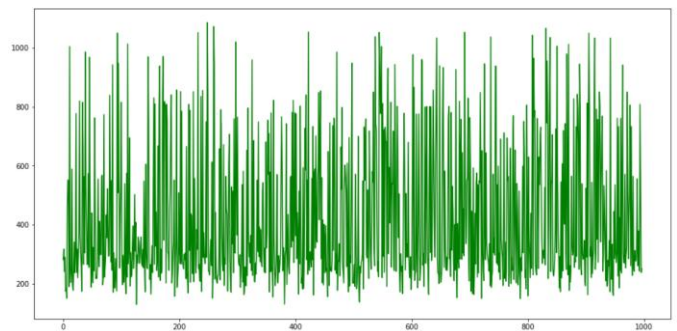
taken.

	Date	Open	High	Low	Close	Volume	Name
0	2006-01-03	211.47	218.05	209.32	217.83	13137450	GOOGL
1	2006-01-04	222.17	224.70	220.09	222.84	15292353	GOOGL
2	2006-01-05	223.22	226.00	220.97	225.85	10815661	GOOGL
3	2006-01-06	228.66	235.49	226.85	233.06	17759521	GOOGL
4	2006-01-09	233.44	236.94	230.70	233.68	12795837	GOOGL

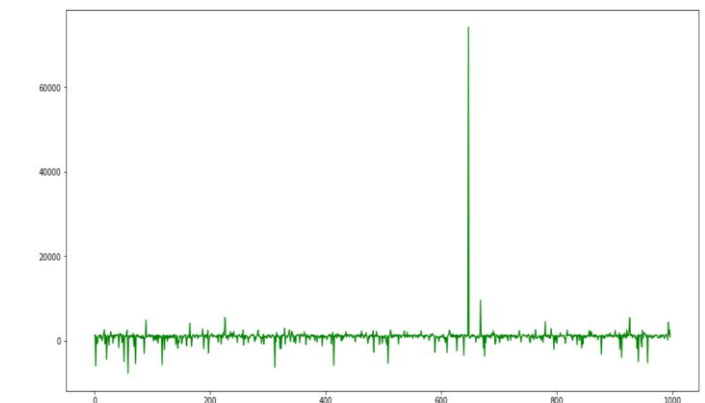
This is a picture of the data present in our csv file. This is the result of using data.head(). It returns the first five rows. Five is the default number of rows returned by data.head() in python pandas library.

	Date	Open	High	Low	Close	Volume
0	2006-01-03	211.47	218.05	209.32	217.83	13137450
1	2006-01-04	222.17	224.70	220.09	222.84	15292353
2	2006-01-05	223.22	226.00	220.97	225.85	10815661
3	2006-01-06	228.66	235.49	226.85	233.06	17759521
4	2006-01-09	233.44	236.94	230.70	233.68	12795837

This is the data after dropping the name from the xlsx file using data.drop('NAME', axis=1).



This graph is telling about the data range.



This graph is representing the actual versus predicted data.

REFERENCES

1. Ashish Sharma, Dinesh Bhuriya, Upendra Singh. "Survey of Stock Market Prediction Using Machine Learning Approach", ICECA, 2017.
2. Pei-Yuan Zhou, Keith C.C. Chan, Member, IEEE, and Carol Xiaojuan Ou, "Corporate Communication Network and Stock Price Movements: Insights From Data Mining", IEEE 2018.
3. Sachin Sampat Patil, Prof. Kailash Patidar, Asst. Prof. Megha Jain, "A Survey on Stock Market Prediction Using SVM", IJCTET 2016.
4. G. Miller, "Social scientists wade into the tweet stream," Science, vol. 333, no. 6051, pp. 1814–1815, 2011.

5. Ryo Akita, Akira Yoshihara, Takashi Matsubara, and Kuniaki Uehara. Deep learning for stock prediction using numerical and textual information. In Computer and Information Science (ICIS), 2016 IEEE/ACIS 15th International Conference on, pages 1–6, 2016.
6. B. Li, K. C. C. Chan, and C. Ou, “Public sentiment analysis in Twitter data for prediction of a company’s stock price movements,” in Proc. IEEE Int. Conf. e-Bus. Eng. (ICEBE), Nov. 2014, pp. 232–239.
7. Lipo W., Shekhar G.: “Neural Networks and Wavelet De-Noising for Stock Trading and Prediction”, Springer, Time Series Analysis, Modeling and Applications Intelligent Systems Reference Library Volume 47, pp 229-247, 2013.
8. Gupta, A. : “Stock market prediction using Hidden Markov Models”, IEEE Engineering and Systems (SCES), 2012 Students Conference on, pp.1-4, 2012. M. Young, The Technical Writers Handbook. Mill Valley, CA: University Science, 1989.

AUTHORS PROFILE



Mr. Arun Kumar, Asst. Professor, Computer Science and Engineering, SRM Institute of Science and Technology, Ramapuram, Chennai, India.

His area of teaching include Cloud computing, Network security, Theory of computation, Compiler Design, Java programming. His research interest include Cloud computing, Network security, Privacy preservation. He has done his b-tech from Anna University in 2006 and later on completed his master of technology from SRM Institute of Science and Technology Ramapuram (formerly known as SRM University), 2010. He is currently teaching in SRM Institute of Science and Technology.



Abhishek Jha, B.Tech Student, Computer Science and Engineering, SRM Institute of Science and Technology, Ramapuram, Chennai, India. He is currently pursuing his b-tech from SRM Institute of Science and Technology, Ramapuram, Chennai, India.

His area of interest include Theory of computation, C language, C++ language, Java programming. He is currently learning web development and machine learning. He knows Hindi, English languages. He likes playing cricket, badminton, football. His technical skills include HTML, CSS, C++, Java, Python language. He has done his internship in Kaashiv infotech, Chennai.



Shashank Shekhar, B.Tech Student, Computer Science & Engineering, SRM Institute of Science and Technology, Ramapuram, Chennai, India. He is currently pursuing his b-tech from SRM Institute of Science and Technology, Ramapuram, Chennai, India.

His area of interest include Theory of computation, C language, C++ language, Java programming. He has also done some projects related to web development. He worked with a local company in Chennai called Kaashiv infotech recently. He likes to read books related to AI and machine learning. He is also interested in cyber security and ethical hacking and is working on his hacking skills recently.



Ashutosh Kumar Singh, B.Tech Student, Computer Science & Engineering, SRM Institute of Science and Technology, Ramapuram, Chennai, India. He is currently pursuing his b-tech from SRM Institute of Science and Technology, Ramapuram, Chennai, India. His area of interest include Theory of computation, C language, C++ language, Java programming. He has a key interest in web

development using MongoDB and Scala. He has also completed his internship in Kaashiv infotech on web development. He also has a key interest in ethical hacking. He likes to read books on administrative subject and one day hopes to work in the administrative area.