

Stock Price Prediction Using Data Mining Techniques

Pavan S, S Usha, Rakshith S, Vijay Joshi, Ravindra Acharya G R

Abstract: Using the past stock knowledge, the paper describes about the development of two models to form short-run predictions for a stock value. The models were refined by the influence of information system index. Advanced mathematical techniques were not able to formulate these models. Investors will use these models to get suggestions and pointers. To check these models we tend to compare the predictions with actual performance of many stocks and obtained trustworthy results. In an exceedingly amount wherever the market went five-hitter down our model yielded a gain of 4.35%.

Index terms: Linear Regression, Support Vector Machine (machine learning), WPI (Wholesale Price Index).

1. INTRODUCTION

The course of concentration is on WPI i.e., arithmetic with computational fixation. The real passing task (MQP), the paper completed an undertaken task to anticipate the costs of stocks.

Chipping away at this undertaking enabled to utilize a large portion of the material we have experienced in our significant field. Strategies utilized in math displaying and numerical examinations classes were useful in approximating the suitable information. Insights classes helped us to comprehend and assess this information. This venture manages propelled math displaying devices and addresses an extremely mind boggling and well known region of the monetary segment. Since our profession objectives are to go into the money related part and get associated with the securities exchange, we think this venture helped us to push ahead mentally. It furnished us with a solid scholastic base and a positive encounter in regards to the securities exchange. In the market a financial specialist can exchange with stocks, choices and prospects. Choice is an agreement that sets a value that you can either purchase or sell a specific stock at an ensuing time. Future is an agreement to sell or purchase a product sometime in the future, at a cost settled upon ahead of time.

Revised Manuscript Received on December 22, 2018.

Pavan S, RRCE, Bangalore

S Usha, Professor and head, CSE, RRCE, Bangalore

Rakshith S, RRCE, Bangalore

Vijay Joshi, Student, Bachelor of Engineering in Computer Science & Engineering in RajaRajeswari College of Engineering, Bangalore

Ravindra Acharya G R, Student, Bachelor of Engineering in Computer Science & Engineering in RajaRajeswari College of Engineering, Bangalore

The thing that matters is that for a fates contract a financial specialist will undoubtedly sell or purchase the ware, while the choices contract gives you the decision to exchange. With a vast piece of the general public attempting to foresee the stock costs, the market is extremely prominent in the present day. They do as such so as to either ensure a monetarily sheltered retirement, procure a living or beat the market. A solid stock portfolio will help accomplish these objectives. With this task we made a numerical model that predicts the cost of offers. This model can enable anybody to select possibly effective stocks and make a solid portfolio. Despite the fact that it is difficult to foresee the future with a 100% assurance, this complex scientific model ought to accomplish a dimension of accuracy worthy by financial specialists and intermediaries alike. Beating the market implies that you are really beating another person. Another person needs to lose with the end goal for you to win. This another person can be an individual simply like you or it very well may be a huge monetary association. These money related organizations have numerous experts and a lot bigger cash-flow to contribute. Achieving this task from the money related and scientific point of view will

help speculators not to lose to the market in this manner levelling the fields. We trust our stock determining models will be helpful for individual financial specialists and retirees searching for a steady future who have no entrance to definite data about the execution of the organizations behind the stocks the corporate and worldwide legislative issues, for instance, the drop in oil costs because of military/political emergency with Russia. Activities all in all are constrained to examining the money related and scientific piece of the financial exchange, which is the reason it is practically difficult to make a 100% exact expectation. We are endeavoring to make the most precise portrayal of what's in store for explicit stocks. We are endeavoring to produce exact conjectures from fifteen to thirty business days. This undertaking will comprise 3-4 organizes so as to make a refined model. The main phase of making this model will be this MQP.

WPI is a region well known in the general public and we additionally have an individual

Published By:

Blue Eyes Intelligence Engineering & Sciences Publication



intrigue. The innovation segment, in any case, is difficult to speak to by demonstrating few stocks. That is the reason we begun taking a shot at the venture by picking ten stocks from the Internet Information Providers Industry of the innovation part. There were a couple of limitations set from the earliest starting point for choosing the stocks. To start with, the organization stocks should be generally steady, accordingly the majority of the stocks we picked have opened up to the world for quite a while now. The stock value information for these organizations is accessible for each working day for as long as year. Second, the value scope of the stocks over five dollars and underneath

hundred dollars. We picked the ten stocks by taking a gander at their costs and placing them in 3 divisions. A value scope of 5 – 20, 20 – 50 or more 50. These stocks with their stay away from little market capitalization stocks with a cost of under five dollars at first on the grounds that their cost can be changed effectively by a financial specialist with an expansive capital. We chose to maintain a strategic distance from vast market top stocks with a cost of in excess of 100 dollars since they are not an exact portrayal of the business, for instance, Google is viewed as a market creator and can't speak to a particular industry of a segment.

Many areas such as review websites, online retail uses sentimental analysis and it is a versatile social media analytics used at the recent times.

Sentiment analysis plays an important role in customer service, management of brand reputation and business intelligence. Sentimental Analysis also plays a vital role in politics as well. Twitter analysis has been used in the Demonetization step taken by the Indian government also sentimental analysis has been helpful for the American 2011 election. Earlier days the prediction of the stock was made by the external experts in a traditional way. Efficient market hypothesis were used instead of prediction algorithms. According to it the prices of the stocks is decided by the new as well as the different patterns. This made a way to develop the prediction algorithms which uses the past data to determine the future stock price which would help the entire stock market to correct the stocks knowing the future price of it. Many researchers have rejected the use of the prediction algorithms and tried to get the patterns of the stock market's behavior.

The different algorithms used in the prediction of the stock price are Linear Regression, Artificial Neural Networks, Naivebayes classifier. Based upon the textual analysis from the twitter feeds, several attempts were done to get a accurate and reliable prediction value.

Another paper proposed by Si, J., Mukherjee, A., Liu, B., Li, Q., Li, H., & Deng, X. expressed methodology to make stock prediction using the sentiments of Twitter Corpus which was very useful for Market Scenario prediction.

This paper used a continuous Dirichlet Process Mixture Model to learn the daily dataset. These both datasets were then put into linear regression to predict the expected values as result.

II. LITERATURE REVIEW

Similarly, Gidofalvi, G., & Elkan, C. proposed another paper with similar ideology which extracted dataset from financial news articles based on which prediction of the Short-Term Movement of the stock market scenario. The

authors classified the movement of stock market scenario into 3 different classes, which are “up”, “down” and “unchanged”. This paper used Naïve Bayesian Text classifier to predict the direction of the movement of stock market scenario, by deriving the key word set from the text

data that is available on various online financial news portals.

III. METHODOLOGY

The dataset comprises of information from the timetable 2000 to 2010 which is utilized for the forecast of Dow Jones Industrial Average Index. The description of this data has been recovered from three distinct sources-

1. News information from Reddit World News Channel. Only the best features of e are considered.
2. The stock information for Dow Jones Industrial Average (DJIA) is considered from the year 2010 to 2017. Yahoo account is utilized for the stock information.
3. Data from the Guardian's peaceful news API.

The main features for multi day are organized in succession alongside a 'name' segment holding '0' if these features lead to the decrease or absence of progress in the stock cost the following day and holding '1' if these features lead to increment in the stock value the following day. The information recovered from the best 25 features and the name is contrasted with get a precise expectation. The content in the preparation dataset is changed over into numeric qualities utilizing 'Sack of Words' model. This model essentially tallies the quantity of events of each word in the given information. The estimations of number events of the words in the information are utilized to shape highlight vectors got from the model.

By using the essential 'Sack of Words' model, which describes the development of the highlight vectors just dependent on the recurrence of words ignoring the request which it processes. This way, we can demonstrate skip gram or a n-gram show which uses and stores the created components.

The n-gram mentions briefly or indirectly to the n words in the given request. This proposed paper utilizes a unigram(n=2) and a bigram(n=2) display which utilizes the tally of sets of two words all together to be stored. The vectors got from these strategies are utilized to prepare the AI models.

Irregular Forest classifier is an administered arrangement

Published By:

Blue Eyes Intelligence Engineering & Sciences Publication



calculation. It makes numerous choice trees dependent on irregular subsamples from the information, each equipped for creating an outcome when given qualities for expectation. The larger the quantity of trees, the larger will be the exactness and lesser is the danger of over fitting contrasted with different models.

A. Random Forest

The Random Forest Classifier produces 'n' number of trees as frail classifiers and unions every one of the trees into a backwoods. At the point when the Random Forest Classifiers is utilized for relapse the mean of coming about qualities from all the choice trees is the subsequent expectation esteem and when it is utilized for arrangement, the subsequent class is the method of the subsequent classes from the choice tree. To characterize another item each tree gives a characterization that can be portrayed as a vote and the class with the most noteworthy votes is picked as the class of the new article.

Random Forest Classifier is a group technique utilized for arrangement or relapse. Irregular Forest Classifier works utilizing a colossal accumulation of de-connected choice trees. In this, the preparation information shapes a lattice as input. Using this framework, a substantial number of new grid with arbitrary components are created. Using every one of these network, a comparing choice tree is shaped for order of the testing information. At the point when the testing information is input, all these choice trees order the info test information and anticipate the class to which the information belongs. The result is discovered dependent on the expectation result which has the most extreme consider the aftereffect of the classifiers. To make expectations, when the preparation is done, the Random of forecasts from all person In AI, support vector machine is a critical model. It is an administered learning model utilized for arrangement and regression [8]. In this model, we are given a lot of preparing examples in which every last one of them is set apart to have a place one of the two categories [8].

B. Support Vector Machine

A Support Vector Machine show speaks to a point in space which is mapped to such an extent that the two unique classes with their components are isolated by as much separation as possible. The picture pixels or the contribution under test is mapped into this space and forecasts are made dependent on the class or classification to which the test input has a place [2]. The graphical portrayal of working of a SVM

$$\text{margin} = \arg_{S_{x,D}} \min(|x \cdot w + b| / \sqrt{\sum_{i=1}^n w_i^2})$$

The separation between the two parallel hyperplanes picked to distinguish the most extreme edge hyperplane is given by, $2/||w||$

The objective of executing Support Vector Machine is to effectively characterize the data. Therefore, The given information focuses are restrained from falling in the edge.

$$y_i = +1 \text{ when } w \cdot x_i + b \geq +1$$

$y_i = -1 \text{ when } w \cdot x_i + b \leq -1$ Hence from the over two disparities,

$$y_i (w \cdot x_i + b) - 1 \geq 0$$

A Quadratic Programming definition of Support Vector machine having hard edge can be given by, $\min w, b (||w||/2)$, with the end goal that $y_i (w \cdot x_i + b) \geq 1 \forall i \in \{1, 2, 3, \dots, n\}$

C. Linear Regression

Linear Regression is basically the earliest algorithm in Machine Learning which every data scientist is familiar with. It is a simple model but everyone needs to master it as it lays the foundation for other machine learning algorithms. It is a very powerful technique and can be used to understand the factors that influence profitability. It can be used to forecast sales in the coming months by analyzing the sales data for previous months. It can also be used to gain various insights about customer behaviour. By the end of the blog we will build a model which looks like the below picture i.e, determine a line which best fits the data. It is a very powerful technique and can be used to understand the factors that influence profitability. It can be used to forecast sales in the coming months by analyzing the sales data for previous months. It can also be used to gain various insights about customer behaviour. By the end of the blog we will build a model which looks like the below picture i.e., determine a line which best fits the data.

D. Ridge Regression

In edge relapse, the cost capacity is modified by adding a punishment identical to square of the

extent of the coefficients. This is equal to stating limiting the cost capacity in condition 1.2 under the condition as underneath

Supplement 1: Constrain on Ridge relapse coefficients So edge relapse puts limitation on the coefficients (w).. How about we see a model utilizing Boston house information and underneath is the code I used to portray direct relapse as a restricting instance of Ridge relapse.

E. Lasso Regression

The cost capacity for Lasso (least supreme shrinkage and determination administrator) relapse can be composed as Supplement 2: Lasso relapse coefficients; subject to comparative compel as Ridge, appeared. Much the same as Ridge relapse cost work, for $\lambda = 0$, the condition above lessens to condition 1.2. The main contrast is as opposed to taking the square of the coefficients, sizes are considered. This sort of regularization (L1) can prompt zero coefficients for example a portion of the highlights are totally ignored for the assessment of yield. So Lasso relapse helps in lessening over-fitting as well as it can help us in highlight determination. Much the same as Ridge relapse the regularization parameter (λ) can be controlled and we will see the impact beneath utilizing malignant growth informational collection in sklearn. Reason I am utilizing disease information rather than Boston house information that I have utilized previously, is, malignancy informational index have 30 highlights contrasted with just 13 highlights of Boston house information. So highlight choice utilizing Lasso relapse can be delineated well by changing the regularization parameter.

The general system fused in Random Forest model to build its security and exactness is called Bootstrap Aggregation which is otherwise called packing. In this method, when given preparing contribution alongside the reactions, arbitrary subsamples are picked persistently from the information exhibited to us with substitution to fabricate choice trees.

Given preparing set- $X = x_1, x_2, x_3, \dots, x_n$.

Given Responses in respect to the above preparing set- $Y = y_1, y_2, y_3, \dots, y_n$.

for i in 1 to n

1. Choose Subsamples from the preparation set with substitution from both X and Y calling these, X_i, Y_i .

2. Create a choice tree F_i for relapse or order utilizing X_i and Y_i .

After the making of various choice trees the expectation for concealed example of information, x' can be made by estimation of the normal of forecast esteems from every one of the trees. What's more, if there should be an occurrence of characterization, by taking the method of anticipated classes from every one of the trees.

IV.CONCLUSION

In this undertaking we made two models for transient expectations of stock costs. We looked at the changed models and abridged the outcomes. To test gauges of our models continuously, virtual \$100,000 was utilized as an interest in various stocks. In a time of 18 business days, we had the capacity to post an addition of over \$4,000 while list overall dove over 5%. In light of this execution we are certain that these models merit the thought of any person who should need to put resources into the financial exchange.

At last, after the intensive investigation of the information acquired, we have outlined extra contemplation for any individual who may choose to get on the venture and keep working a similar way.

REFERENCES

- Breiman, L. (1996). Bagging predictors. *Machine Learning*, 24(2), 123-140.
- Burges, C.J. (1998). A tutorial on Support Vector Machines for pattern recognition. *Data mining and knowledge discovery*, 2(2), 121-167.
- Hearst, M. A., Dumais, S.T., Osuna, E., Platt, J., & Scholkopf, B. (1998). Support Vector Machines. *IEEE Intelligent Systems and their applications*, 13(4), 18-28.
- Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.
- Du, W., & Zhan, Z. (2002, December). Building decision tree classifier on private data. In *proceedings of the IEEE international conference on Privacy, security and data mining-Volume 14* (pp. 1-8). Australian Computer Society, Inc.
- Liaw, A., & Wiener, M. (2002). Classification and regression by RandomForest. *R news*, 2(3), 18-22.
- Scholkopf, B., & Smola, A.J. (2002). *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press.
- Hsu, C. W., Chang, C. C., & Lin, C.J. (2003). *A practical guide to support vector classification*.
- Thissen, U., Van Brakel, R., De Weijer, A. P., Melssen, W. J., & Buydens, L. M. C. (2003). Using support vector machines for time series prediction. *Chemometrics and intelligent laboratory systems*, 69(1), 35-49.

AUTHORS PROFILE



Retrieval Number:F12910476S519/19©BEIESP

S. Usha, is working as a professor and head, CSE, RRCE. Graduated from Manonmaniam Sundaranar University, in CSE during the year 1998. She obtained her Master degree in CSE and PhD degree from sathyabama university in the area of Mobile Ad Hoc Networks in the year 2013. She has 54 publications in International and National conferences, 22 publication in national journal and international journals in the area of Mobile Ad hoc Networks and wireless security. Most of the publications are having impact factor cited in google scholar (h index and i10index), Microsoft etc.. Received fund from AICTE under NCP scheme and SERB(DST). Received best teacher award from Lions club in the year 2010&2012. Developed Centre of Excellence lab in IoT with industry collaboration. Organized many conferences, FDPs and Technical Talks. Associated with ISTE, CSI, IEEE, IAENG, IDEAS and IACSIT. Reviewed papers in IJCs and CiiTjournals. Acted as a TPC member in MIRA'14 IoTBDS '17 and IoTBDS'18 Portugal. Chaired sessions in FCS'14, ICISC'13 & ICCCT'15, ICCCT'17 and IoTBDS'18.



10. **Pavan S**, pursuing Bachelor of Engineering in Computer Science & Engineering in RajaRajeswari College of Engineering. An Active member of Computer Society of India. Awarded with Student Project of the Year 2017-18 under IEAE and is indexed in IEAE Digital Library of Academic Projects (DLAP).



11. **Ravindra Acharya G R**, pursuing Bachelor of Engineering in Computer Science & Engineering in RajaRajeswari College of Engineering. An Active member of Computer Society of India and active member of IEEE.

Stock Price Prediction using Data Mining Techniques