

STOCK MARKET PREDICTION USING SENTIMENT AND TIME SERIES ANALYSIS IN R

Submitted to: E.A.S. (UB)

December 2021

Group No. 2

Abstract

This paper studies the possibilities of forecasting stock market prices of firms using the sentiments captured via web scrapping. We have experimented with stock market price of Tesla and Moderna using sentiment analysis and ARIMA model. An accuracy analysis was also carried out with a R-squared value of each of the model to evaluate how each of them fared in the forecasting. The aim is to help reduce participants in loss while investing using the twitter data. The stock data was pulled of the Yahoo Finance API. The sentiments were obtained off the sentences of tweets from twitter. Results obtained has proved that the ARIMA model has good R-Squared value for short term prediction.

Keywords

R, Stock Market Prediction, Time Series Analysis, Sentiment Analysis, ARIMA Model

1. Introduction

The aim of this paper is to forecast the price of the Tesla and Moderna stock by using time series analysis by using the ARIMA model. It uses the advance computing of computers as a tool to extract data from the web and then parse it on tokens for further computation. The parsed text is further analyzed for sentiments thereby giving each a numeric value ranging from -1 to 1. The penultimate goal is to forecast the stock price of the previously mentioned firms by using the tweets and news data sentiments with the already existing live stock price. For this purpose, we have used historical data of opening price, closing price and highest volume traded for Tesla and Moderna. The stock price is affected by various factors like current trade scenarios, people's liking, the company's performance, pandemics, etc.

The rest of this paper is organized as follows: Section 2 explains some basic concepts. Then, Section 3 explains the system architecture. Followed by Section 4 that explains the ARIMA model. Finally, Section 5 which delivers the results and conclusions of the work.

2. Background

2.1 Basic Concepts

- **Data Mining –**
To find anomalies, patterns, correlations within large datasets to predict future outcomes.
- **Time Series Analysis –**
Time series analysis is a specific way of analyzing a sequence of data points collected over an interval of time.
- **Sentiment Analysis –**
Sentiment analysis is contextual mining of text which identifies and extracts subjective information in source material and helping a business to understand the social sentiment of their brand, product or service while monitoring online conversations.
- **Supervised Learning –**
It is training a computer on a specific input data that has been labelled for a specific output.
- **Web Scrapping –**
It is used for data harvesting. Also known as the process of using bots to extract data from a website.
- **Stock Market –**
It refers to public markets that exist for buying, issuing and selling stocks that trade on a stock exchange or over the counter. To raise capital of a firm.

3. System Architecture

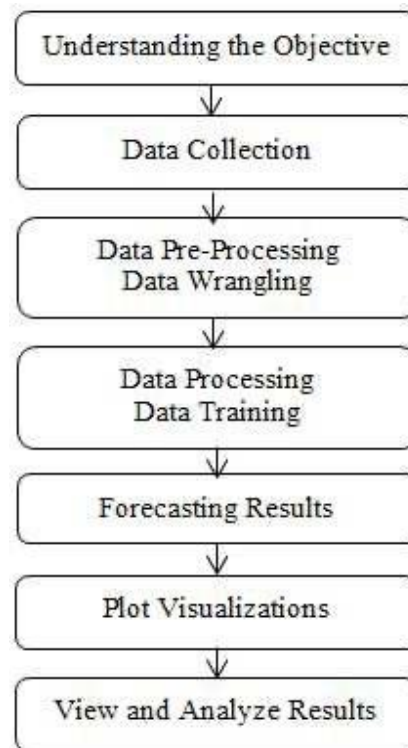


Fig. 3. System Architecture

- **Understanding the Objective –**
The aim of our project is to build a model that predicts the stock price of a given firm while it extracts words from twitter on related topics. Our goal is neither to make billions off the system nor waste billions too. But the objective is to help stock market investors by giving them a direction in taking a decision or not. Whether to buy/hold/sell a stock by providing the result in terms of visualizations.
- **Data Collection –**
This is the process where we used a python script to scrape data off twitter. A sample of the python script used to scrape data is shown below. Also, we used Yahoo finance to get the stock price data for the corresponding interval of time.

```
1 from Sweet.scweet import scrape
2 from Sweet.user import get_user_information, get_users_following, get_users_followers
3
4 data = scrape(words=['amc'], since="2021-01-15", until="2021-08-09", from_account = 'reuters', interval=1, headless=False, display_type="Top", save_images=False, lan
5 | resume=False, filter_replies=False, proximity=False)
```

Fig. 3.1 Web Scrapping Script (Python)

- **Data Pre-Processing –**

This is the stage where the acquired data is processed into final datasets to work on. Cleaning the dataset is the focus in this stage. The overview of the dataset is shown below. There are many columns out of which we used the “Embedded text” column as the main feature for sentiment analysis as that contained the tweets of people. The data pre-processing was also used on the stock price data to make it ready and be combined into a value vector.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
	UserScreen	UserName	Timestamp	Text	Embedded_text	Emojis	Comments	Likes	Retweets	Image link	Tweet URL						
1	HumanLife /	@HumanLife	2020-06-01T	Human Life	WATCH: Tech billionaire and			1	3	7	https://pbs.twimg.com/media/1267573018626318342						
2	TESLARATI	@Teslarati	2020-06-01T	TESLARATI	Tesla broadens Model 3 shop with All-Weather Protection Kit and			3	25	241	https://twitter.com/Teslarati/status/1267588091260305408						
3	Dan Trent	@trent_dan	2020-06-01T	Dan Trent	Credit due, Tesla caught the Germans napping. Based on first go in			3	3	72	https://pbs.twimg.com/media/1267570521610338309						
4	UWish	@jwiah	2020-06-01T	UWish	Join us this morning (in about 2 hours) for a chat about #EV driving	👉👉👉		2	5	5	https://pbs.twimg.com/media/1267570408979223552						
5	San Francisco	@sfchronicle	2020-06-01T	San	After Tesla got permission from local authorities to reopen its			10	39	115	https://pbs.twimg.com/media/1267579801292279816						
6	Adam	@cipil429	2020-06-01T	Adam	In letter to state, Tesla says it.Ads below required number of			5	15	53	https://twitter.com/cipil429/status/1267572281175310336						
7	Fabrizio Busti	@FabrizioBusti	2020-06-01T	Fabrizio	Jay Leno drove the			5	5	10	https://pbs.twimg.com/media/126759721111755776						
8	Knee Of The	@KneeOfThe	2020-06-01T	Knee Of The	Here.Ads my take on #Cybertruck on Jay Leno.Ads Garage. it.Ads			1	5	15	https://pbs.twimg.com/media/126754050067984256						
9	Chris Grey	@3rdwaven	2020-06-01T	Chris Grey	@Tesla			39	17	83	https://twitter.com/3rdwaven/status/1267564535218491395						
10	FSD Pilot	@jchylow	2020-06-01T	FSD Pilot	Anyone notice that AutoPilot has improved drastically on winding	👉👉👉		10	5	55	https://twitter.com/jchylow/status/1267571419413708800						
11	ACKIO,AC	@Kristennet	2020-06-01T	ACKIO,AC	Hey if you see the tell him we appreciate the fight he puts in for	👉👉👉		5	5	60	https://twitter.com/Kristennet/status/1267602482965786624						
12	GigaCam (id	@GigaCam1	2020-06-01T	GigaCam	Just applied for the	👉👉👉		3	19	19	https://twitter.com/GigaCam1/status/1267576668566785						
13	MikelnVegas	@fhdogs	2020-06-01T	MikelnVegas	Hey			1	6	6	https://pbs.twimg.com/media/1267561108639936512						
14	Eva CyberFo	@EvaFoxi	2020-06-01T	Eva	Tesla Sales In Germany Could Benefit With The Possible Extra	👉👉👉		2	5	26	https://pbs.twimg.com/media/1267541888040648704						
15	eden	@vlexwier	2020-06-01T	eden	IVE MISSED YOU	👉👉👉		1	4	5	https://twitter.com/vlexwier/status/1267547633477812225						
16	Back Roads	@jamwest	2020-06-01T	Back Roads	Why, we please L&B have some more of this SR+ efficiency. 75			1	4	4	https://pbs.twimg.com/media/1267560132809912320						
17	Tesmanian.c	@Tesmanian.c	2020-06-01T	Tesmanian.c	Tesla TSLA Market Cap Surpassed All German Legacy Automakers			1	10	75	https://pbs.twimg.com/media/1267586991366092480						
18	Discover EV	@discover_e	2020-06-01T	Discover EV	Europe overtakes China on EV investment! https://bit.ly/3xveeCt			2	1	2	https://pbs.twimg.com/media/1267551828566302722						
19	Nick van Raa	@MisterNid	2020-06-01T	Nick van	Made a little youtube video on why people should buy an electric			2	2	2	https://twitter.com/MisterNickS/status/1267549715010334720						
20	Electrek.Co	@ElectrekCo	2020-06-01T	Electrek.Co	Rivian Adventure Network: electric pickup maker hires Tesla staff			1	25	89	https://pbs.twimg.com/media/1267601978743566343						
21	Trinotor 202	@EanMich	2020-06-01T	Trinotor	Replying to			1	1	4	https://twitter.com/EanMich/status/1267582523836413440						
22	Watt.Ads Ga	@wattsgara	2020-06-01T	Watt.Ads	Wow just heard the	👉👉👉		1	1	1	https://twitter.com/wattsgara/status/1267594769554146306						
23	Elon.Ads Bra	@ElonBrain	2020-06-01T	Elon.Ads	Replying to	👉👉👉		1	1	1	https://pbs.twimg.com/media/1267562876262424848						
24	Center for Ai	@Ctr4Auto	2020-06-01T	Center for	Do we know if this			3	9	11	https://pbs.twimg.com/media/1267595548611264512						
25	Nikki Fried S	@RealNikkiF	2020-06-01T	Nikki Fried	Replying to	👉👉👉		2	12	103	https://pbs.twimg.com/media/1267595831558967299						
26	Tesmanian.c	@Tesmanian.c	2020-06-01T	Tesmanian.c	Tesla Sales In Germany Could Benefit With The Possible Extra			2	12	103	https://pbs.twimg.com/media/1267546723645194240						
27	Jayam DVD	@jydeh9	2020-06-01T	Jayam				3	1	3	https://twitter.com/jydeh9/status/12675339294794499						
28	Dr. Americus	@amreed2	2020-06-01T	Dr.	The			3	1	3	https://twitter.com/amreed2/status/1267587150238765057						
29	Frida Kahlo	@FridaKahlo	2020-06-01T	Frida Kahlo	Elon our BIG little country will have no limits soon. Catalonia will			1	1	4	https://twitter.com/FridaKahloPDC/status/1267558292735373313						
30	Ross Gerber	@GerberKav	2020-06-01T	Ross Gerber	Replying to			1	5	5	https://twitter.com/GerberKawasaki/status/126757019157743616						
31	Carl Koomet	@Ckoomet	2020-06-01T	Carl Koomet	Replying to	👉👉👉			4	4	https://twitter.com/Ckoomet/status/1267564276031664129						
32	LaneM2000	@LaneM2000	2020-06-01T	LaneM2000	Day 52 of tweeting				2	2	https://twitter.com/LaneM2000/status/1267553352503021825						
33	Doviz	@doviz	2020-06-01T	Doviz					1	1	https://twitter.com/doviz/status/126757202304981962724						

Fig. 3.2 Tweets data (TESLA)

- **Data Processing –**

To process the data, we use the ARIMA(p,d,q) model. Generally, stock investors use the auto regressive and moving average models to forecast the future trends. Highlights here would be estimation, forecasting and identification. These steps are repeated recursively until an optimal model is identified for prediction. R language provides auto.arima() method to forecast the time series data according to ARIMA(p,d,q) model.

- **Forecasting Results –**

The process of predicting the future by relying upon the past and current data is called as forecasting. Various prediction techniques are used by stock analysts to predict the future stock trends value. The ‘forecast’ package offered in R was used to predict the future trends which took in values off the sentiment score and past historical stock price data. The “SAS” package of R was used for sentiment analysis thereby giving each sentence a sentiment score. All of these were fed into the ARIMA model which then forecasted the results for time series predictions. It also offered exponential smoothing and space models.

- **View and Analyze Results –**

This is done to evaluate and view the outcome of the model. Screenshots of the evaluation and results are provided further in the section. Investors can view the results and graphs with a comparison view and then invest in the stock. They can use this as an “assistance” to buy/sell/hold a particular stock.

4. ARIMA MODEL

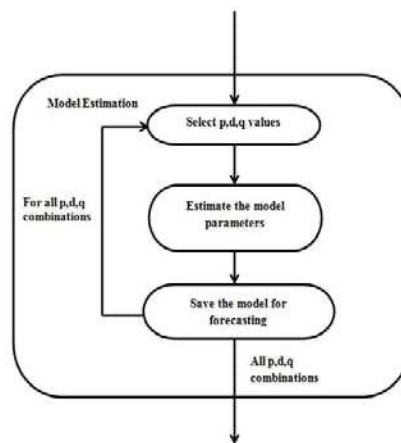


Fig. 4 Model Estimation of ARIMA

- In ARIMA model, the identification is to be accomplished using auto co-relation function and partial auto co-relation function in order to identify p, d and q standards. For any realistic time, sequence generally p, d and q values vary between 0 and 2, but model estimation is executed for all probable combinations of p, d and q values.
- ARIMA() Function in R –
Arima() function automates the inclusion of a constant. By default, d = 0 or the value of d = 1. A constant will be included if it improved the AICc value; for d > 1 the constant is always omitted. If allowdrift= FALSE is specified, then the constant is only allowed when d = 0.

In ARIMA model, the future value of a variable is a linear combination of past values and past errors, expressed as follows:

$$Y_t = \phi_0 + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \varepsilon_t - \theta_1 \varepsilon_{t-1} - \theta_2 \varepsilon_{t-2} - \dots - \theta_q \varepsilon_{t-q}$$

Fig. 4 Arima function Math

Where, Y_t is the actual value and ε_t is random error at t , Φ_1 and θ_j are the coefficients, p and q are integers which are often referred to as autoregressive and moving average resp.

5. Results and Conclusion

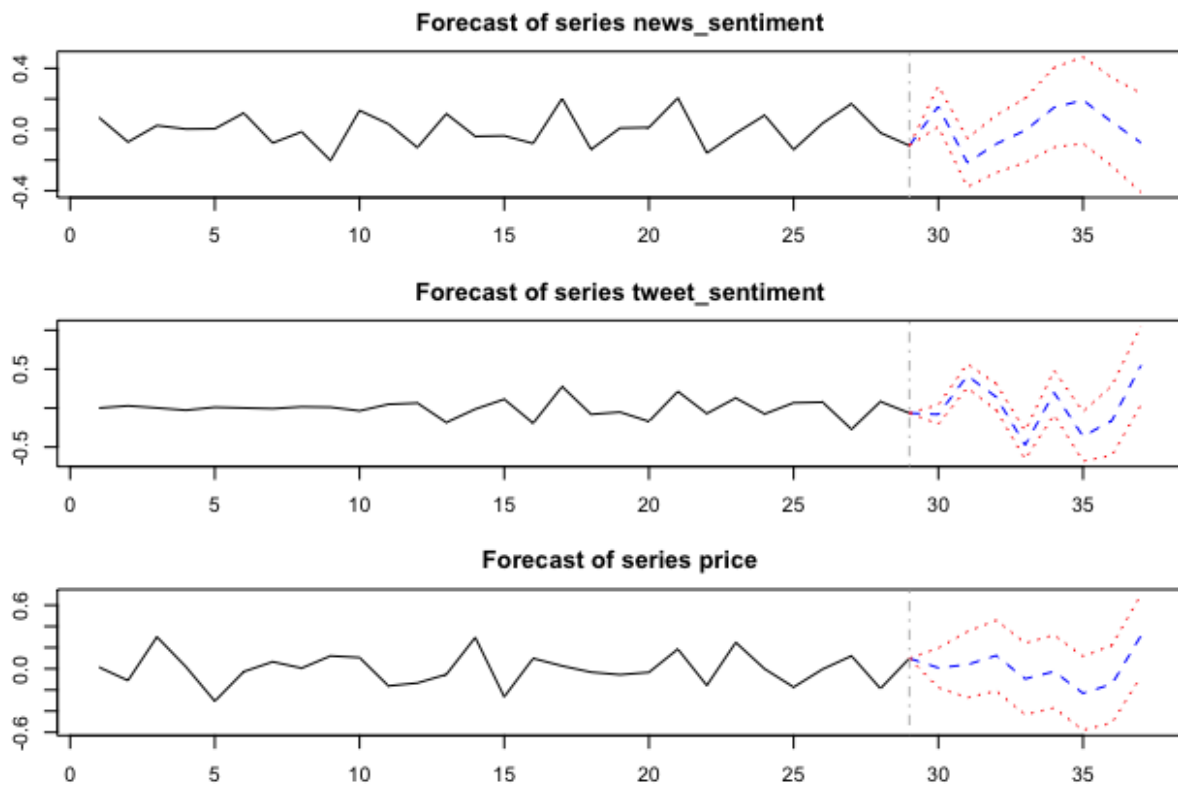


Fig. 5.1 Forecasting News and tweet sentiments and the price

The above figure shows us how the news and tweets of users for the word “Tesla” will be forecasted continuing the current and previous trends.

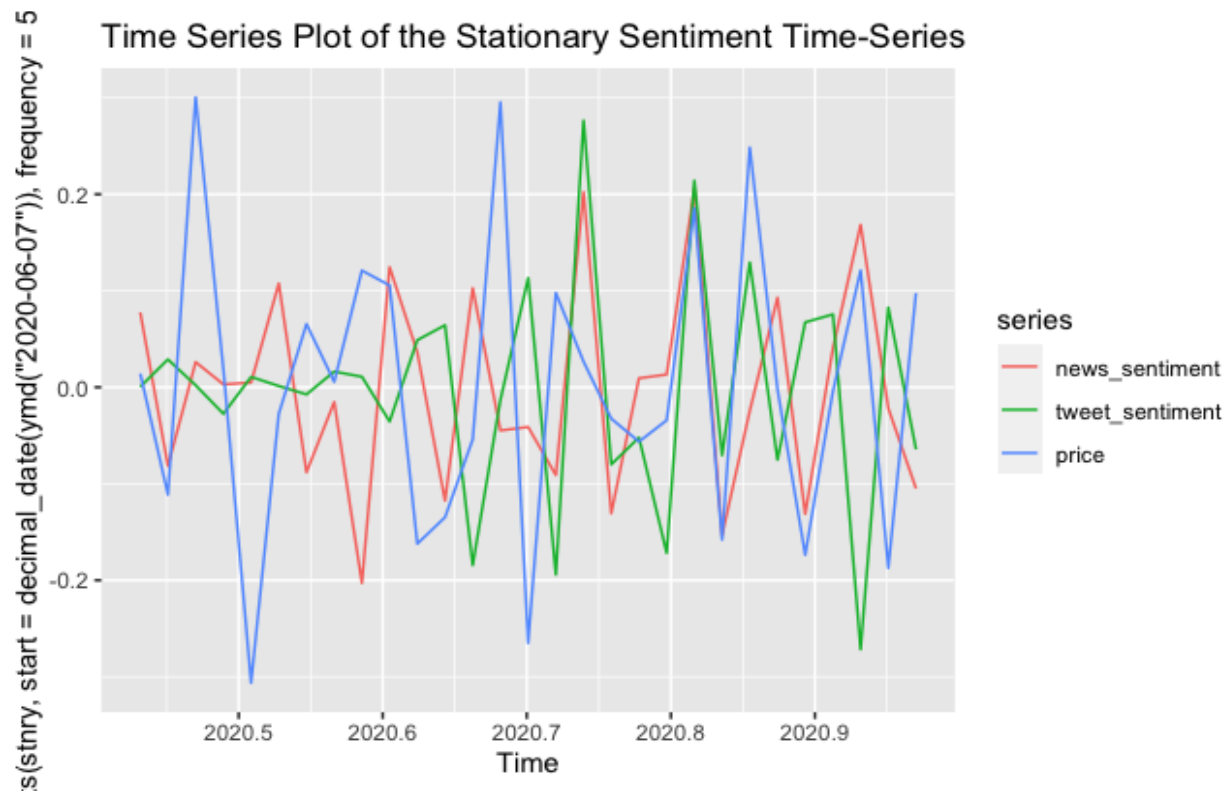


Fig. 5.3 Graphing the tweets and price (TESLA)

Plotting the sentiment score w.r.t. the price shows us the correlation of how the stock price of TESLA has been greatly affected by the tweets.

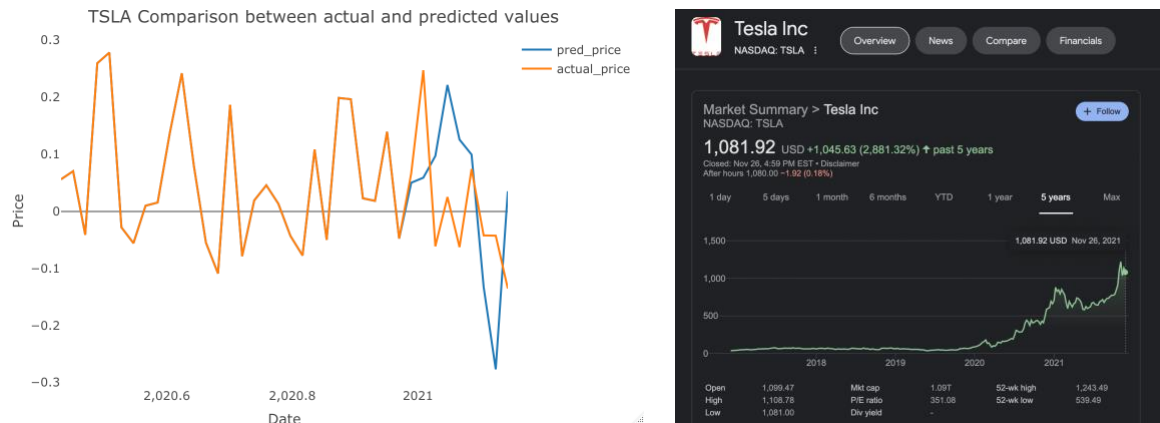


Fig. 5.4 Prediction Comparison of our model (LHS) to real time stock price TESLA (RHS)

As we can see that on the (LHS), the blue line indicates the predicted price of TESLA's stock by our model is close to the real time stock price of TESLA (RHS). Post 2021 start, the model was given to predict the stock and it did perfectly.

R-Squared Value – 54.81 (This determines how accurate the model is)

Conclusion –

- In this paper we attempted to predict the stock market for TESLA and Moderna using the historical stock price data and encompassing the impact of sentiments using tweets and news headlines of the same. The experimental results obtained demonstrated the potential of the ARIMA model in short term prediction. This could guide the investors in investing wisely on whether buy/sell/hold that stock.
- The factors affecting the R-Squared value are regional investing trends, pandemics, stock's current value, etc.
- Investors can take the risk of investment easily when they have an idea about the future value of the stocks.

6. References

- [1] Henrique, B. M., Sobreiro, V. A., & Kimura, H. J. E. S. w. A. (2019). Literature review: Machine learning techniques applied to financial market prediction. pp. 120, 226-251.
- [2] Qasem A. Al-radaideh, Adel Abu Asaf, Eman Alnagi, "Predicting stock prices using data mining techniques", The International Arab Conference on Information Technology 2013
- [3] Li Bing, Chan, K. C. C., C. Ou, "Public sentiment analysis in Twitter data for prediction of a company's stock price movements", 11th IEEE International Conference on e-Business Engineering (ICEBE), November 2014, pp. 219-239.
- [4] Han, J., Kamber, M., Jian P., "Data mining concepts and techniques". San Francisco, CA: Morgan Kaufmann Publishers, 2011.
- [5] A. J. Conejo, M. A. Plazas, R. Espnola and B. Molina. "Day- ahead electricity price forecasting using the wavelet transform and ARIMA models", IEEE Transactions on Power Systems, 2005, pp. 1030–1042
- [6] Banerjee, D., "Forecasting of Indian stock market using time- series ARIMA model", 2nd IEEE International Conference on Business and Information Management (ICBIM), January 2014, pp. 110-135.
- [7] Schumacher, R.P., Chen, H.: Textual analysis of stock market prediction using breaking financial news: the AZFin text system. ACM Trans. Inf. Syst. 27, pp. 1–51 (2009)
- [8] Ayodele A. Adebisi, Aderemi O. Adewumi, Charles K. Ayo, "Stock price prediction using the ARIMA model", 16th IEEE International Conference on Computer Modelling and Simulation (UKSim), March 2014, pp. 100 -112.
- [9] Huang, W., Nakamori, Y., Wang, S.-Y. J. C., & research, o. (2005). Forecasting stock market movement direction with support vector machine. pp. 32(10), 2513-2522
- [10] A. Meyler, G. Kenny and T. Quinn, "Forecasting Irish Inflation using ARIMA Models", Central Bank of Ireland Research Department, Technical Paper, 3/RT/1998.
- [11] P. Pai and C. Lin, "A hybrid ARIMA and support vector machines model in stock price prediction", Omega vol.33 pp. 495-510, 2005
- [12] Fama, E.F.: Random walks in sotck market prices. Financ. Anal. J. 51(1), pp. 74–80 (1995)
- [13] Chen, W., Zhang, Y., Yeo, C. K., Lau, C. T., & Lee, B. S. (2017). Stock market prediction using neural network through news on online social networks. Paper presented at the 2017 international smart cities conference (ISC2).