

# 2020sc04239.pdf

*by*

---

**Submission date:** 20-Mar-2023 03:44PM (UTC+0530)

**Submission ID:** 2041577378

**File name:** 2020sc04239.pdf (1.65M)

**Word count:** 4868

**Character count:** 25175

**A**

**Report On**

**Stock gain forecasting for day opening using ML and NLP  
Dissertation**

**DSECLZG628T: Dissertation**

by

Sarvsav Sharma

2020SC04239

**1**  
**Dissertation work carried out at**

**Tata Consultancy Services Pvt. Ltd.**



**BIRLA INSTITUTE OF TECHNOLOGY & SCIENCE PILANI (RAJASTHAN)**

A Report  
On  
Stock gain forecasting for day opening using ML and NLP Dissertation

1 Submitted in partial fulfillment of the requirements of the M.Tech Data Science and Engineering programme

By

Sarvsav Sharma  
2020SC04239

Under the supervision of

Amit Sharma – Assistant Consultant  
Tata Consultancy Services



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE  
Pilani (Rajasthan) INDIA

## ACKNOWLEDGEMENTS

This project has been possible because of many individuals' consistent assistance and support. I would like to express my gratitude to them for all of their support.

Primary thanks to my supervisor, **Amit Sharma**, for his guidance throughout this project. Our meetings were energy filled and productive.

**Vinaya Sathyanarayana**, for giving me continuous feedback throughout the project stages, and helped me to understand the business value of the project.

I value my friends' helpfulness and upbeat attitudes; their motivation and support have been essential to my achievement.

<sup>6</sup>  
Finally, I would like to offer my sincere gratitude to the Divine and my cherished family for their unfailing support and inspiration while I completed this endeavor.

**BIRLA INSTITUTE OF TECHNOLOGY & SCIENCE, PILANI FIRST  
SEMESTER 2022-23  
DSECLZG628T DISSERTATION**

<b>Dissertation Title:</b>	Stock gain forecasting for day opening using ML and NLP by suggesting to buy or sell the quantities
<b>Name of Supervisor:</b>	Amit Sharma
<b>Name of Student:</b>	Sarvsav Sharma
<b>ID No. of Student:</b>	2020SC04239

**Key Words:**

Stocks gain, loss recovery, machine learning for stocks, mean squared error, easy trading for individuals, stock grouping for gains

**Abstract:**

The prediction of stock market price trends is a constantly popular topic for academics in both the financial and technical professions since it is one of the most crucial areas of concentration for traders to generate rapid money.

With an emphasis on short-term price trend forecasting in the intraday market, the aim <sup>6</sup> of this study is to construct a cutting-edge prediction model for price trend forecasting. Traders can decide how much stock they can trade for the least amount of loss based on the results.

The present methods for identifying the optimal movement produce poor stock predictions because they either model numerical data or statistical data. Consequently, it's important to take into account both elements in order to get a more accurate outcome.

A machine learning model that estimates the appropriate quantity of stocks to purchase or sell depending on trader risk level in order to achieve maximum gain or minimal loss is trained to handle this issue. In order to make predictions with a higher degree of accuracy, this algorithm scrapes real-time news and shares price information from the internet.

Due to the project's interactive flow, even a person with little computer experience can readily use it. The information for the news and stock prices is gathered using python requests and nsetools from reputable online sources like economonictimes. By putting the program on the cloud, the scraping may be done continuously.

A dataset with 1M entries would be an optimal size for a training model. The column names have been updated, filtered, cleared of outliers and null values, and scaled appropriately for a better result. The sentiment analyzer Vader is used to convert the textual data since it makes it possible to determine not just the positive and negative score but also the compound score and neutral score.

To gather data on model improvement, feature selection uses bar graphs, pair plots, line plots, and heat maps.

Several ML algorithms are applied to the dataset, and even the results are calculated using the auto h2o model that is capable of finding the best model. The model with highest accuracy is the SVR poly for the current dataset.

The best performing model is selected and saved as the news prediction model.pkl file for deployment based on the performance of multiple models. The stored model is also utilized in the Python project to forecast the optimal number of stocks that a trader should purchase or sell in order to realize the most profits.

The project stands out because it takes factors like platform fees and risk level and trains on both textual and numerical data in order to provide greater insights for traders regarding market gains.



(Signature of Student)

Name: Sarvsav Sharma

Date: 05/03/2023

Place: Gurgaon



(Signature of Supervisor)

Name: Amit Sharma

Date: 05/03/2023

Place: New Delhi

**Problem statement:**

Stock gain forecasting for day opening using ML and NLP by suggesting to buy or sell the quantities

**Completeness criteria:**

Project should be able to predict the quantity of stocks that a trader can buy or sell with maximum gain and less risk factors.

**Best performance:**

Model is performing better than the previous model in our organization and it performs better than the h2o model that compares various results and chooses the best model among them.

**Success criteria:**

Model is able to successfully predict the outcome for buying or selling the quantity of stocks for maximum gains.

**Challenges Faced:**

1. Collection of headlines data has some characters outside of ascii text and needs to remove them manually.
2. Handling of sentiment or news score for the weekend data, as there are no movements in stock prices because of the market closed.
3. Timeout issue with the python nsetools, and an active machine that runs every minute to collect the data. A cron job is suitable to run on the cloud with minimal resources and is best to collect the real time data for the stocks.
4. Data web scraping, or creating a universal web scraping code that works for all years for economictimes website to collect the latest headline, as they have different templates for the five years old pages.

**BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE, PILANI**

**CERTIFICATE**

This is to certify that the Dissertation entitled **Stock gain forecasting for day opening using ML and NLP** and submitted by Mr **Sarvsav Sharma** having IDNo. **2020SC04239** in partial fulfillment of the requirements of Dissertation, embodies the work done by him under my supervision.



Signature of the Supervisor

Amit Sharma

Assistant Consultant

Tata Consultancy Services Pvt. Ltd.

Place: Gurgaon  
Date: 05/03/2023

## 1 **Table of Contents**

ACKNOWLEDGEMENTS	2
ABSTRACT	3
CERTIFICATE	6
<b>CHAPTER 1: INTRODUCTION</b>	10
<b>CHAPTER 2: SCOPE OF PROJECT</b>	12
<b>CHAPTER 3: LITERATURE SURVEY</b>	13
<b>CHAPTER 4: DATA COLLECTION</b>	14
<b>CHAPTER 5 : MACHINE LEARNING MODEL PERFORMANCE</b>	23
<b>CHAPTER 6 INTEGRATING WITH PYTHON PROJECT</b>	27
<b>CHAPTER 7 : CONCLUSION AND INFERENCE</b>	29
<b>CHAPTER 8: REFERENCES</b>	30
<b>APPENDIX</b>	31

## **List of Figures**

Figure 1 : Approach for calculating gains	10
Figure 2 : Sample output from python nsetools	13
Figure 3 : Nifty50 data for every minute for the last 5 years and showing all the features	14
Figure 4 : News archive from economic times	14
Figure 5 :Top headlines for the day	15
Figure 6 :Extracted headlines from the website to parse in HTML format	15
Figure 7 : Distribution plots for the features	18
Figure 8 :Scatter plots for the features	19
Figure 9 :Box plot to determine the outliers	20
Figure 10 : Heat Map	20
Figure 11: Observation from various models	21
Figure 12 : Experimental Output	24
Figure 13: Saving the model	27
Figure 14: Running within python project	28

## **List of Tables**

Table 1 : Factors influencing stock price	12
Table 2 : Collected Dataset	15
Table 3 : Clean and Prepared Dataset	16

## Chapter 1: Introduction

The value of stocks always fluctuates dramatically over time, making stock market prediction an intriguing study subject. Investors do two different forms of study before buying a stock.

The fundamental analysis is one of the first approaches and very common. In order to determine whether to invest or not, investors consider factors such as the intrinsic worth of the stocks, the state of the market and economy, the political environment, etc. On the other side, technical analysis evaluates equities by looking at data produced by market activity, such as previous prices and volumes. Typically of attempting to determine a security's fundamental worth, technical analysts instead utilize stock charts to spot patterns and trends that could predict how a stock will act in the future.

Many methods for predicting stock movements have been developed throughout the years. Initially, stock trend predictions were made using traditional regression techniques. Non-linear machine learning methods have also been applied since stock data may be characterized as non-stationary time series data.

With a forget gate present, the LINEAR, POLY is similar to a long short-term memory (SVM), however it has fewer parameters than the SVM since it lacks an output gate. The vanishing gradient issue that arises when using a conventional scaler is addressed with LINEAR, POLY. The time sequence is erratic and disordered. Most of the forecasting model that uncovers the complex connection between financial information about an industry and its stock price is beneficial. The financial news in addition to the existing records concerning the firm is used to forecast future stock prices.

Semantic and linguistic traits may be extracted using a variety of ways. The following are a few of them: OpinionFinder, SentiWordNet, Linguistic Inquiry and Word Count (LIWC), Google Profile of Mood States (GPOMS), R sentiment analysis, and Python NLP package. In this approach, the sentimental score is also calculated based on news headlines, in addition to the statistical data for the model to produce more reliable results.

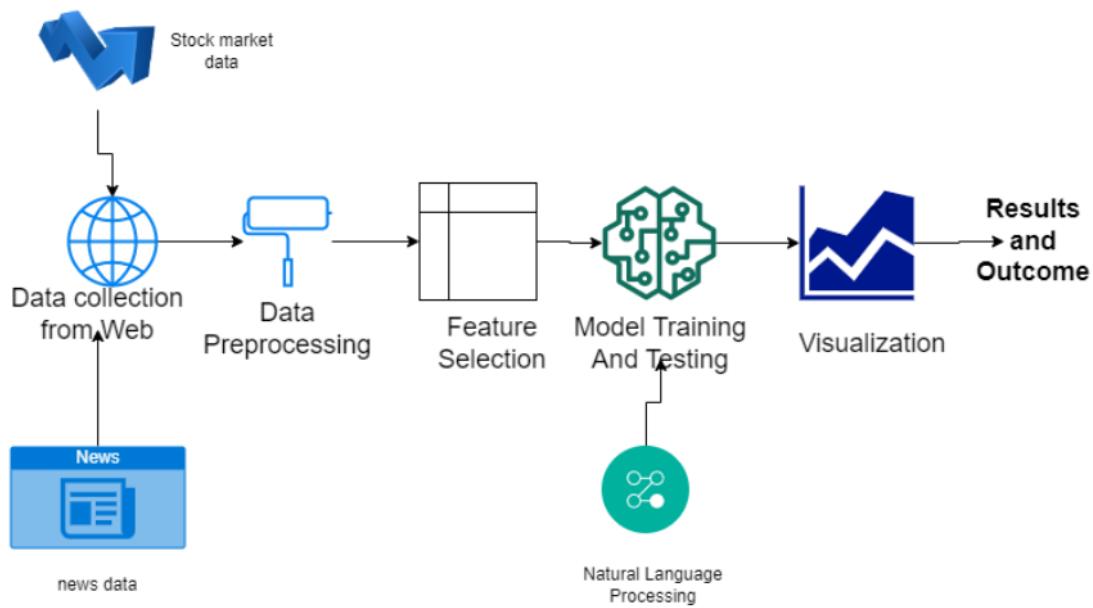


Figure 1: Approach for calculating gains in stock market using machine learning and natural language processing

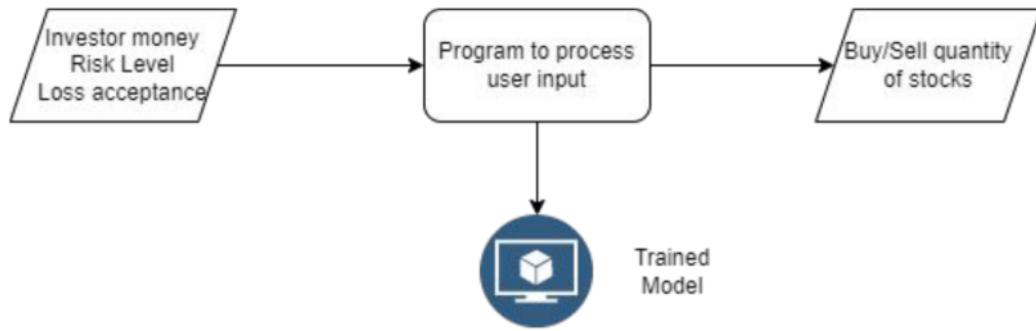


Figure 2: Data flow diagram for the project

## **Chapter 2: Scope of project**

Goal of the project is to help the traders to buy or sell the correct quantity of shares to obtain maximum gains for a day and recover their losses. As a human, whenever we notice a small increase in price of stock from the buying rate, we sell it and do not capitalize the maximum gain that we can have from the share at the end of the day. And, the same is valid for selling also, as we keep believing that the stock will grow up and at the end of day it results in a loss.

The project can be used by individual traders who are investing money and small scale trading companies to help their customers in wise trading.

There are many models available on the internet that predict the value for the upcoming stock, however they are difficult to use and, on the other hand, don't guide you when not to invest.

The existing models either predict whether the stock price will go up or down but don't advise you on the number of stocks that you can buy for minimum loss and maximum profit.

The project is open for anyone and recommended to use for the new traders who prefer to trade wisely with minimal risks.

The model has trained for a very limited number of stocks or popular companies like HDFC bank. To be used on a huge scale, the model needs to be trained for a bigger set of companies. It is limited in scope of predicting the results for the set of trained companies.

The inputs for the project are the name of the stock in which the trader would prefer to invest, the risk level on a scale from one to five considered higher, the number is higher the risk, the platform fees, and the acceptance loss. Once the inputs are provided, the program will decide whether the user needs to buy, sell, or quit for today based on the news and other statistical parameters. The output will be the quantity of stocks that a user can buy to have maximum gain at current time and price.

The output result is very understandable as it just prints the quantity of stocks to buy, sell, or not to invest based on the user risk level.

The small scale companies helping the traders can use it to track the status for multiple stocks from different companies on which the model has been already trained, and the individual traders can safely focus on one company or stock to trade for maximum profit.

## Chapter 3: Literature Survey

Many investing decisions are influenced by financial markets worldwide. As a result of the geopolitical, social, and economic developments occurring throughout the world, stock markets experience significant shifts over time.

The majority of studies employed quantitative information, such as history or current prices, as predictor variables to forecast the current stock price. The utilization of the massive amounts of unstructured textual data created by the web in the form of news stories published on websites, user opinions on social media, and blogs written by professionals in the field of financial investments received less attention.

The investors and the financial institutions are the stakeholders and they in turn cause financial risks in investments. As a result, academics began examining the causal connection between different market conditions and the related changes in stock values.

Several machine learning models were compared while forecasting the price direction of companies and indices on the Indian stock market in a work by Patel et al. [7]. They were naive Bayes, random forest, support vector machine, and artificial neural network. Ten technical metrics based on open, high, low, and close prices were utilized as input data, expressed as continuous values between -1 and 1. The results showed that random forest performed better than the other three models overall, with an accuracy rate of 83.56%, while naive-Bayes performed the worst, with an accuracy rate of 73.3 %.

Now, there are different approaches to investing in stock like buy, sell, and hold. Nelson, Pereira, and Oliveira [8] investigated the effectiveness of LSTM networks for forecasting changes in stock price. By analyzing performance using the criteria accuracy, precision, recall, and F-measure, the outcomes were compared to other investment methods and machine learning models. Multi-layer perceptron, random forest, and a pseudo-random model were the models selected for comparison. Additionally, they evaluated the success of buy-and-hold strategies and an optimistic strategy in which stocks were acquired if their prices had increased in the previous time step and sold in the next time step. The period of time selected was 15 minutes. The researchers' suggested LSTM model surpassed the baseline comparisons, averaging 55.9% accuracy on average.

## **Chapter 4: Data Collection**

Information gathering is the project's most vital stage. Nsetools is only one of the numerous Python modules available for real-time data collection. To get headline data, you may scrape websites like Economic Times. The information about data gathering is provided in this section. While historical stock data may be retrieved from the internet, real-time news headline data is collected using a web scraper from the website [economictimes.com](http://economictimes.com).

In both the instance of the market data and the news data, the data is quite information-rich. There are several elements that might affect the stock value and may not be covered by the news, therefore it is not advisable to rely just on one news source. And the top six factors mentioned below that have the potential to influence stock prices

- Supply and demand
- business metrics
- details about the industry and advances
- broader market trends
- Geographical context
- The economy's interest rates
- Investor mindset

Table 1: Factors influencing stock price

The most popular Python web scraping packages or frameworks are Beautiful Soup, Scrappy, and requests. requests mainly utilized for websites with java script. As a consequence, data extraction is carried out using a Python script.

### **a. Data Extraction**

The stock website's data is extracted using a Python script and nsetools, and the result is a CSV file. Additionally, [economictimes.com](http://economictimes.com) is utilized to read archived information as a web page for news data, and the beautifulsoup library of Python is used to read the HTML information. The top news stories for the day were then obtained using the regex. The script has been adjusted to operate with the website's most recent ten years' worth of data.

The features and their specifications are covered in this document. It is now essential to choose the characteristics from the stock market data and important phrases like vaccination

that have an influence on the stock market. Due to the output obtained from the website's collection of several needed variables and the possibility of trivial information, only the following feature is taken into account.

- Date - date of observation
- Open - opening value for the index on that date
- Close - closing value for the index on that date
- High - highest value achieved by the index on that date
- Low - lowest value by the index on that date

This information is gathered throughout the last 5 years and that can help to find the gains for nifty50 stock. Data with all the features manually collected for the news and script is ready to collect the data from nse (national stock exchange).

```
>>> q = nse.get_quote('nifty50') # it's ok to use both upper or lower case for codes.
>>> from pprint import pprint # just for neatness of display
>>> pprint(q)
{'adhocMargin': None,
 'applicableMargin': 12.5,
 'averagePrice': 1999.82,
 'bcEndDate': None,
 'bcStartDate': None,
 'buyPrice1': 1999.45,
 'buyPrice2': 1999.4,
 'buyPrice3': 1999.35,
 'buyPrice4': 1999.15,
 'buyPrice5': 1999.1,
 'buyQuantity1': 50.0,
 'buyQuantity2': 209.0,
 'buyQuantity3': 22.0,
 'buyQuantity4': 1.0,
 'buyQuantity5': 24.0,
 'change': 25.35,
 'closePrice': None,
 'cm_adj_high_dt': '01-DEC-14',
 'cm_adj_low_dt': '30-MAY-14',
 'cm_ffm': 190659.16,
 'companyName': 'Infosys Limited',
 'css_status_desc': 'Listed',

'dayHigh': 2010.0,
'dayLow': 1972.0,
'deliveryQuantity': 258080.0,
'deliveryToTradedQuantity': 51.54,
'exDate': '02-DEC-14',
'extremeLossMargin': 5.0,
'faceValue': 5.0,
'high52': 2201.1,
'indexVar': None,
'totalTradedValue': 22914.16,
'totalTradedVolume': 1145811.0,
'verMargin': 7.5}|
```

>>>

A10	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	
179141	2017-02-01	8761.5	8761.55	8759.4	8759.7	0	8760.68	8759.8	8759.943	8760.36	8760.483	8760.182	8760.012	8759.865	8762.538	8760.68	8758.822	8759.714	8759.528	8759.927	8760.169	
179142	2017-02-01	8760.95	8762.25	8760.7	8761.45	0	8760.71	8759.92	8759.84	8760.585	8760.639	8760.322	8760.129	8759.968	8762.579	8760.71	8758.841	8759.679	8759.554	8759.952	8760.170	
179143	2017-02-01	8762.65	8763.45	8759.5	8760.8	0	8760.87	8760.355	8759.897	8760.745	8761.309	8760.745	8760.444	8760.223	8763.186	8760.87	8759.554	8759.808	8759.877	8760	8760.229	
179144	2017-02-01	8764.4	8764.6	8763.3	8762.85	0	8761.85	8760.96	8760.207	8760.862	8762.339	8761.41	8760.939	8760.621	8765.01	8761.85	8758.69	8760.096	8760.502	8760.059	8760.33	
179145	2017-02-01	8764.3	8764.55	8763.4	8764.4	0	8762.76	8761.495	8760.613	8760.762	8762.993	8761.935	8761.359	8760.972	8765.58	8762.76	8759.94	8760.44	8760.963	8760.111	8760.42	
179146	2017-02-01	8763.7	8764.6	8762.7	8764.6	0	8763.2	8761.94	8760.793	8760.757	8763.229	8762.256	8761.651	8761.231	8765.771	8763.2	8760.629	8760.747	8761.201	8760.133	8760.496	
179147	2017-02-01	8761.85	8763.5	8760.7	8763.8	0	8763.38	8762.045	8761.073	8760.725	8762.769	8761.282	8761.678	8761.29	8765.352	8763.38	8761.403	8761.022	8761.29	8760.145	8760.511	
179148	2017-02-01	8761.5	8762.15	8761.15	8761.9	0	8763.15	8761.021	8761.287	8760.77	8762.346	8762.058	8761.659	8761.313	8765.616	8761.35	8760.684	8761.265	8761.12	8760.152	8760.523	
179149	2017-02-01	8762.5	8763.25	8760.4	8761.45	0	8762.77	8762.31	8761.568	8760.847	8762.397	8762.139	8761.76	8761.424	8762.77	8760.429	8761.28	8760.202	8760.562	8760.562		
179150	2017-02-01	8759.5	8762.75	8759.5	8762.35	0	8761.81	8762.285	8761.6	8760.912	8761.432	8761.659	8761.477	8761.24	8764.564	8761.81	8758.69	8761.355	8761.28	8760.195	8760.355	
179151	2017-02-01	8761.2	8761.4	8759.5	8759.5	0	8761.31	8762.255	8761.73	8761.027	8761.354	8761.575	8761.443	8761.237	8763.316	8761.21	8759.304	8761.572	8761.179	8760.211	8760.569	
179152	2017-02-01	8760.35	8761.45	8760.5	8761.25	0	8761.01	8762.195	8761.7	8761.057	8761.02	8761.533	8761.306	8761.152	8763.052	8761.01	8758.96	8761.606	8761.271	8760.212	8760.566	
179153	2017-02-01	8760.55	8761.15	8760.05	8760.35	0	8760.82	8761.855	8761.613	8761.701	8761.059	8761.207	8761.212	8761.059	8762.821	8760.82	8758.819	8761.649	8760.217	8760.566	8760.566	
179154	2017-02-01	8761.5	8761.55	8760.5	8760.6	0	8760.62	8761.695	8761.747	8761.327	8761.075	8761.26	8761.248	8761.133	8762.018	8760.62	8759.222	8761.068	8762.36	8760.246	8760.573	
179155	2017-02-01	8760.75	8762.3	8760.5	8761.45	0	8760.87	8761.34	8761.814	8761.471	8761.097	8761.185	8761.097	8761.715	8760.87	8760.025	8761.595	8761.229	8760.252	8760.576		
179156	2017-02-01	8760.75	8761.15	8759.8	8760.55	0	8760.78	8761.045	8761.769	8761.492	8760.895	8761.091	8761.131	8761.064	8761.559	8760.78	8760.001	8761.558	8761.03	8760.256	8760.576	
179157	2017-02-01	8760.7	8760.7	8759.85	8760.55	0	8760.85	8760.93	8761.973	8761.477	8761.487	8760.873	8760.762	8760.077	8761.029	8760.129	8761.516	8760.85	8760.184	8761.561	8760.26	8760.579
179158	2017-02-01	8760.55	8761.15	8760.5	8760.35	0	8760.85	8760.835	8761.607	8761.422	8760.736	8760.935	8761.011	8760.984	8761.516	8760.85	8760.184	8761.585	8761.182	8760.263	8760.576	
179159	2017-02-01	8759.5	8760.5	8759.2	8760.9	0	8760.45	8760.355	8761.28	8761.422	8760.324	8760.674	8760.822	8760.642	8761.411	8760.455	8759.489	8761.375	8761.058	8760.26	8760.574	
179160	2017-02-01	8759.5	8759.75	8755.1	8759.55	0	8759.48	8760.175	8760.72	8761.21	8758.73	8759.806	8760.270	8760.372	8763.175	8759.48	8755.785	8761.406	8760.801	8760.133	8760.518	
179161	2017-02-01	8758.7	8760.45	8756.1	8756.1	0	8759.07	8759.525	8760.387	8761.09	8758.8	8759.605	8760.019	8760.212	8762.559	8759.075	8755.581	8761.141	8760.510	8760.176	8760.573	
179162	2017-02-01	8760.95	8761.9	8758.35	8758.35	0	8759.12	8759.855	8760.372	8761.09	8759.516	8759.849	8760.135	8760.283	8762.707	8759.12	8755.533	8760.873	8760.521	8760.111	8760.513	
179163	2017-02-01	8761.2	8762.65	8759.7	8761.75	0	8759.25	8760.65	8760.307	8761.017	8760.708	8760.995	8760.268	8760.37	8763.074	8759.255	8760.442	8760.560	8760.522	8761.022	8760.572	
179164	2017-02-01	8748.2	8764.7	8733.7	8773.7	0	8756.99	8758.72	8759.35	8760.270	8756.118	8759.732	8758.76	8759.211	8766.573	8756.99	8747.407	8760.21	8758.559	8758.683	8760.047	
179165	2017-02-01	8730.05	8748.35	8748.24	8747.6	0	8751.82	8755.65	8759.39	8758.495	8747.429	8758.263	8755.171	8756.434	8775.581	8751.82	8738.059	8759.18	8751.082	8752.704	8756.443	
179166	2017-02-01	8740.4	8744.45	8730.35	8732.25	0	8748.16	8753.615	8756.003	8757.73	8745.086	8750.597	8753.325	8754.907	8772.19	8748.16	8741.234	8759.73	8750.11	8751.671	8755.871	
179167	2017-02-01	8740.2	8740.2	8733.85	8739.25	0	8744.01	8751.955	8754.856	8746.66	8747.454	8748.706	8751.684	8753.506	8764.707	8744.01	8723.313	8756.727	8749.201	8750.775	8755.003	
179168	2017-02-01	8734.5	8742.5	8731.8	8739.4	0	8738.67	8748.96	8752.92	8754.897	8740.472	8746.123	8749.538	8751.696	8730.929	8738.67	8741.811	8755.666	8747.511	8749.209	8753.422	
179169	2017-02-01	8723.65	8734.8	8723.55	8734.8	0	8733.76	8745.375	8750.4	8752.454	8734.864	8742.037	8746.473	8733.76	8721.045	8745.167	8743.965	8745.649	8750.444	8750.444		
179170	2017-02-01	8717.75	8723.55	8715.15	8723.55	0	8731.25	8741.535	8747.517	8750.858	8729.076	8737.756	8742.7	8746.023	8749.607	8731.25	8723.893	8758.09	8741.809	8748.746	8748.746	
179171	2017-02-01	8723.75	8727.75	8718.25	8718.25	0	8728.64	8738.4	8745.29	8749.162	8728.501	8735.737	8741.780	8744.244	8744.606	8728.64	8721.674	8750.228	8738.608	8740.226	8745.179	
179172	2017-02-01	8730.05	8727.3	8727.85	8727.65	0	8727.21	8735.61	8743.447	8747.797	8730.017	8735.232	8739.815	8743.178	8739.675	8727.21	8714.745	8748.227	8738.191	8739.801	8744.567	
179173	2017-02-01	8737.2	8737.8	8731.5	8737.3	0	8727.75	8735.21	8741.89	8746.63	8725.411	8732.411	8735.359	8739.488	8742.609	8741.59	8727.75	8731.845	8747.565	8748.165	8744.25	
179174	2017-02-01	8740.75	8737.65	8733.75	8737	0	8730.95	8732.155	8740.433	8745.473	8734.158	8735.964	8742.136	8745.555	8730.95	8737.545	8745.663	8748.124	8739.593	8743.941	8743.941	
179175	2017-02-01	8740.5	8745.4	8737.2	8737.35	0	8735.15	8738.2	8739.407	8744.422	8736.272	8736.789	8739.414	8741.981	8744.286	8735.15	8725.016	8744.016	8734.811	8739.763	8740.83	

Nifty50 data for every minute for the last 5 years and showing all the features

[BUDGET'23](#) [ETPrime](#) [Markets](#) [News](#) [Industry](#) [Rise](#) [Politics](#) [Wealth](#) [MF](#) [Tech](#) [Jobs](#) [Opinion](#) [NRI](#) [Panache](#) [ET NOW](#) [More](#)

---

**Archives**

2023  
January | February | March | April | May | June | July | August | September | October | November | December

2022  
January | February | March | April | May | June | July | August | September | October | November | December

2021  
January | February | March | April | May | June | July | August | September | October | November | December

2020  
January | February | March | April | May | June | July | August | September | October | November | December

2019  
January | February | March | April | May | June | July | August | September | October | November | December

2018  
January | February | March | April | May | June | July | August | September | October | November | December

2017  
January | February | March | April | May | June | July | August | September | October | November | December

2016  
January | February | March | April | May | June | July | August | September | October | November | December

2015  
January | February | March | April | May | June | July | August | September | October | November | December

2014  
January | February | March | April | May | June | July | August | September | October | November | December

2013  
January | February | March | April | May | June | July | August | September | October | November | December

2012  
January | February | March | April | May | June | July | August | September | October | November | December

2011  
January | February | March | April | May | June | July | August | September | October | November | December

2010  
January | February | March | April | May | June | July | August | September | October | November | December

2009  
January | February | March | April | May | June | July | August | September | October | November | December

---

**MOST READ** **MOST SHARED** **MOST COMMENTED**

For Hindenburg Research, Adani Group's case a man-made disaster in making  


Sukhoi-30, Mirage 2000 fighter planes crash in Morena; Two pilots rescued, 1 killed  


Gautam Adani: Richest man of Asia in the eye of a storm  


LIC doubles down on Adani amid short seller row  


Should investors sell Adani stocks or buy the dip? Sanjiv Bhasin answers  


[More](#)

---

**Not to be Missed**

Polls: Trump plans stops at early-voting states  


2024 H-1B registrations from March 1-7  


Budget 2023: Skilled workforce, clean energy  


Budget 2023: Expectations on  


News archive from economic times



The Economic Times homepage features a top navigation bar with links for Home, BUDGET'23, ETPrime, Markets, News, Industry, Rise, Politics, Wealth, MP, Tech, Jobs, Opinion, NRI, Panache, ET NOW, and More. A search icon is also present. Below the header, there's a main menu with Home, Business News, and Archive. The main content area displays various news stories, including one about Adani Group's land acquisition in Gurugram and another about the Indian Space Research Organisation's Mars mission. A sidebar on the right lists the 'Most Read', 'Most Shared', and 'Most Commented' stories, such as the Hindenburg disaster and Adani's wealth. A 'Not to be Missed' section highlights political news from Assam and Bihar.

## Top headlines for the day

Extracted headlines from the website to parse in HTML format

### **b. Cleaning and preparing data**

The gathered information is then cleansed and formatted in line with the model's requirements. All unnecessary data is eliminated along with duplicate and null values. Furthermore, superfluous columns are eliminated. The most important and labor-intensive stage of machine learning technology.

1

Python's built-in logic and a number of statistical techniques are used to clean and prepare the data. For instance, the price was a character type rather than a number. The most important phase in creating the derived columns from the present columns is finished at this point. As an example, the precision value has been improved to up to 2 decimal places. In the news headlines section, the special symbols have been removed.

Outliers are addressed separately during this step since they might provide unexpected consequences. The proper data type has been allocated to the data since type is the foundation of computer languages. For instance, date has been converted from a string to a datetime type and generic statistics have been determined. Data standardization has also been carried out.

HDFC Ergo raises Rs 350 crore through 10-year NCDs																						
G15	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	
2	1/1/2017	Here are : Café Cour How our Advantage Akhilesh: The tug-of-war between the SF partitions will doubtless continue																				
3	2/1/2017	Another : Soccer-Gi Tech Mah NHL-Wan Maruti, H Apple Inc. Amar Singh Soccer-Isr Dangal of AG VoltE BIP decid UPDATE 3 Grill Pervi Tower cor UPDATE 1 Justice TS Airtel, Vox UPDATE 1 Breaking : Startups I CII wants																				
4	3/1/2017	Seeking v Cabinet in Soccer-Be Will oppo Interest ri Railways i UPDATE 1 Everyone WhatsApp! UPDATE 1 UPDATE 2 Soccer-Fn UPDATE 6 UPDATE 6 UPDATE 2 Enough tr Narendra Man who What nin Soccer-M Sunil Mu																				
5	4/1/2017	Japanese, Kia, Chang TMC's Sui Battle for Will the B Trouble b Today's G Pro-wrest Airtel offy Mynta fr Vodafone Ratan Tat. Samajivat NPCI plan Phonebar Phonebar Taxation i Will tax cr UPDATE 1 I-T dept r UPDATE :																				
6	5/1/2017	Nokia bag Nokia bag FDC Inflow Finance Google CI MakeMy Cash crust Here is w/ Yuukii Dig Googt to Tennis-Mi Tata Powi Impact of India's de Supermar Tax-wary. Prickly iss Economy Novells ta Finance ?																				
7	6/1/2017	UPDATE 2 Supreme After IAS - Alpine Ski DCW nott Journey to Cash crust Cyrus Mis SoftBank Elections. This place Motor rac Small com Tennis-Dj Supreme Edelweiss Edelweiss Note ban Saathi: He El India p Post note																				
8	7/1/2017	Motor rac Trial to rev No truce : McLaren : Government GDP to gr Take note Donald Tr Poke Me: Is in sink or BSP's secr Samsung? How top 1 Kamalapa Russian O amany-sj China was Portugal i Post arre: Vistara Al Congress																				
9	8/1/2017	Buy OneP Why Cedi Poor the I-BJP slams FRI again: AgustaWt Liquor ou Not enro Global wr Mahendri Demonet: Here is all Notes bar How a bu Cycle kitter best Centre to CBI bring: Government Now, han February																				
10	9/1/2017	Battlefield Bombay T Nimbus C Checkout I am still S Government SP family. CBI oppos Black moi Black moe No chang States rec RBI starte Sinking to 2017 will Ness Digit ONGC clo Now, yan' Now, yatr Pravasi B																				
11	10/1/2017	Top corp Anuj Pur. Will India AO Smith Market tc World Barl Basi Filkari to d With 76.2 Stellar Va Doo sled Ajay Pirar Sonalkar Canadian UPDATE 1 CAT 2016 icy winds Mamata f Ronald																				
12	11/1/2017	1 UPDATE 1 WRAPUP WRAPUP Tata Trust BIP suppe The BCCI Independent Lowa Katan Tat. Britain offa Zaha has Akhilesh-Soccer-Isr-Soccer-Be CYCLING - 1FACTBOX- Soccer-Be INTERVIEW! Aircel mar McDonald Forget da																				
13	12/1/2017	Everyone: World is a Centre to Altel Smaller p Hopeful c Hopeful c IDFC Altel Niti Aayog esports! I BPL seeks Bharti Air Congress Indian car Supreme Automata BPLC raise SC dampe Supreme After Suni Congress																				
14	13/1/2017	UPDATE 1 Dad's val Reliance J Virtual re Airtel to s No Jallika VC need VCs need idea to ra DM to h Venture c Telecom Police Bitcoin pr Scrapp all c Price wa Ratan Tat. Sebi may Uttar Prat Tamac Nac Rupee re																				
15	14/1/2017	Soccer-Ni PAC build Rule arm Hero Mot HDFC Erg Barcelona Infosys na Hipcoach Reliance J India can Soccer-Gi Suriname Nations C Newer IIB. With 4 Soccer-Tu Maruti Su No more Aircel or M. Tesla unv																				
16	15/1/2017	Why inde Soccer-Di Guinea Bi UPDATE 1 PM Modi N Chandri Moving to NFL-Raid Barak OT The soon- Rock shov, Jadavpur Books the President Soccer-Se Soccer-Se Soccer-Se Soccer-Gi What is t Buy Garra: Soccer-N																				
17	16/1/2017	Pro Wrest Jubilant BIP target Global sp Warburg : Warburg : Mobile g Services India rest: You aren't May the F Central sc US Reserv WhatsAppi Maharash N Chandri: NPPA man Economy Government US push i																				
18	17/1/2017	After Con Daalchi mri For RIL, d UPDATE 1 Let Jamn Aarti Pdharman Privacy cc BSP's "Mu Utter Prai Alok Vern FSSAI brir Golf-Mcln Phillips no Apple tez it's advan Til岡a ge Soccer-Ay No conce Bad loan Micromax Micromax Telecom 1-IT-Depar Raisina Di China first Ayen pen Utar Prai You may s You may : Free WiFi should be Should be Soccer-M Kumar Vt Demand																				
19	18/1/2017	Tata Mot Staggerin Aion Capi Alon may UPDATE 1 6 BPOs, 2 For Cash- Warburg / Idea Cellu Congress? ILBFS in t IL&PS in t Lenders n CEEA App in an Indi RMZ to bi RMZ to bi Airtel app Black moi Black mo UP polls: Reliance I If you elec Note ban: UPDATE 1 Printing o Airtel ma Akhilesh Salling-Sle Only 20% Universal Growth c Small car Mallikarj Govt to t Havells se Orient Gr HRD Mini SBI may n Frenzied CCI slaps																				
20	19/1/2017	Rs 3.87 cr Managing UPDATE 1 RBI close! Another c CCI seeks Hillary Cli Samajivat Trai calls f Shrikant S Soccer-Ge Bharti Mr Cooperati Cooperati It's a blac: RSS raises How ban M&M bu M&M bu Gerrard ti INTERVIE																				
21	20/1/2017	Here are 1 Zaria Was Pathshar Drop Flp: Why 'don Jabra Hali Here are / Meet this Discover! Maximum Post date How bre in Manpa Uttarakha Can Andhi Resurgen UPDATE 1 Now sma: Golf-Con UPDATE 1 Soccer-Ri																				
22	21/1/2017	Govt sets UPDATE 2 Soccer-M UPDATE 1 UPDATE 2 State-run Parivesh Alpine skii UPDATE 6 Boos for Boost for India may Neither C Monaco g East of crucial M UPDATE 1 Donald Tr CM Harisli Double In Doubte Ir																				
23	22/1/2017	Everyone Liverpool Congress Nice ham FDI outfit Delhi Hig Cyber crir Nolda exp Soccer-Tu Qureshi ji Lawyers v Banks tell Coal scan Donald Tr Budget 2i Supreme Donald Tr US Presid FRBM rev Supreme Tamil Nar																				
24	23/1/2017	24/1/2017 After Con Daalchi mri For RIL, d UPDATE 1 Let Jamn Aarti Pdharman Privacy cc BSP's "Mu Utter Prai Alok Vern FSSAI brir Golf-Mcln Phillips no Apple tez it's advan Til岡a ge Soccer-Ay No conce Bad loan Micromax Micromax Telecom 1-IT-Depar Raisina Di China first Ayen pen Utar Prai You may : You may : Free WiFi should be Should be Soccer-M Kumar Vt Demand																				
25	24/1/2017	HCL to fo BIP silent: BAA-Livid Huawei a Donald Tr Apple por Kalyan Sir NFL-Roef Aircel urg UPDATE 1 UPDATE 1 Uber, Ola Apple exr Real Mad Idea Celu Sanjiv Pur Despite b Apple wa Chief just Hit by Jio, Parthasai																				
26	25/1/2017	26/1/2017 NHAI ris: Do we rei Lessons ir Tata grou PM Nare India, UAI Apple see What list Fast-form USSR sup Bill Gates PM There Starbucks BIP shed: Essar Gro Essar Gro CII fails to Apex cou Bharti Air Donald Tr After Jalli																				
27	26/1/2017	UEA in to Women's don't Why com NFL-Youn? Nokia bag Nokia bag Soccer-Hi Murderec Hallilovic i idea set t: Stepmpt Aam Aadhi Private er UK ag Helion Ve Helion Ve NFL-Supe Indrapras Idea, Vod PM Nare																				
28	27/1/2017	27/1/2017 UAE in to Women's don't Why com NFL-Youn? Nokia bag Nokia bag Soccer-Hi Murderec Hallilovic i idea set t: Stepmpt Aam Aadhi Private er UK ag Helion Ve Helion Ve NFL-Supe Indrapras Idea, Vod PM Nare																				
29	28/1/2017	28/1/2017 BIP no loi With Con. With safe ITC Q3 ne Mallya ha Mukhtar, BSNL evai Angry over UPDATE 1 Southern 2 Will regio When gal: Eccleslon UID debi UID debi With prof few jobs. Oppositio Uber pric IndiaPost																				
30	29/1/2017	29/1/2017 Late goal! Lacklustre UPDATE 1 Budget 2i Trump is i Innovatio Traditions Dining in Brexit blo Does the Why Mou Today, if y Buy Apple East India A food cri A lookbac All you w: A year aft Homemal Soccer-Se Soccer-Se																				

### c. Converting data from textual to numerical form

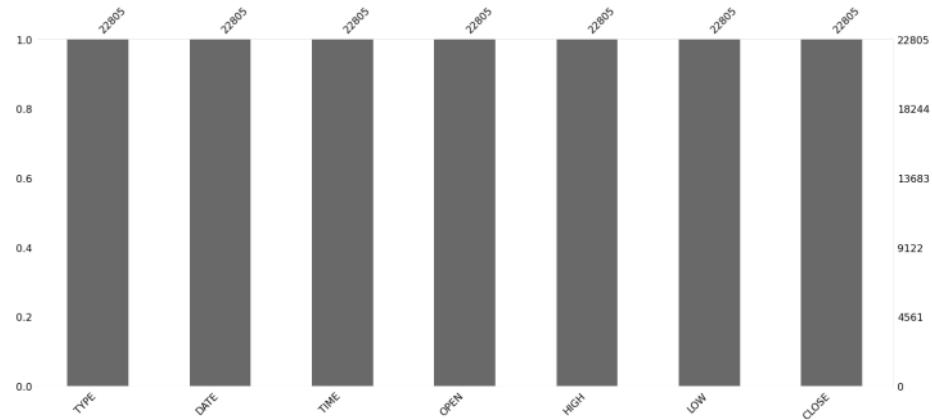
6

The data is in the form of text and to make it ready for the model it needs to be converted into numerical form, as the model operates on numerical data. There are multiple options to convert the data into numerical form. News headlines can have two types of sentiments, either positive or negative. However, Vader[10] helps to calculate the compound and neutral score also. The resultant numerical field can be further used as a numerical input to predict the results. The reason for choosing VADER because with it there is no requirement for general processing the text as it handles everything and abstract a lot of work.

H	I	J	K	L
headlines	compound	negative	neutral	positive
ET Recom	0.9988	0.087	0.797	0.117
It is time v	0.9958	0.096	0.792	0.112
IT hardwa	-0.998	0.106	0.8	0.094
Saudi Ara	0.9918	0.09	0.803	0.107
Aamir Kha	0.998	0.074	0.828	0.098
Indira Gan	-0.9958	0.1	0.81	0.09
If China do	0.9983	0.08	0.813	0.107
Zomato sh	0.9955	0.046	0.887	0.067

A	B	C	D	E	F	G	H	I	J	K	L
	prevclose	open	high	low	last	close	headlines	compound	negative	neutral	positive
1/1/2016	1082.15	1082.4	1090.25	1076.15	1088.7	1088.75	ET Recom	0.9988	0.087	0.797	0.117
1/4/2016	1088.75	1084	1084	1068.1	1068.5	1070.5	It is time v	0.9958	0.096	0.792	0.112
1/5/2016	1070.5	1070.2	1074.8	1061.35	1062	1062.4	IT hardwa	-0.998	0.106	0.8	0.094
1/6/2016	1062.4	1056.65	1076.75	1056.65	1067.55	1067.1	Saudi Aral	0.9918	0.09	0.803	0.107
1/7/2016	1067.1	1060.1	1064.9	1049.7	1052.55	1056.2	Aamir Kha	0.998	0.074	0.828	0.098
1/8/2016	1056.2	1061.95	1064.5	1057.25	1062	1062.35	Indira Gar	-0.9958	0.1	0.81	0.09

1  
Table 2 : Clean and Prepared Dataset



Updating the data to the updated data type

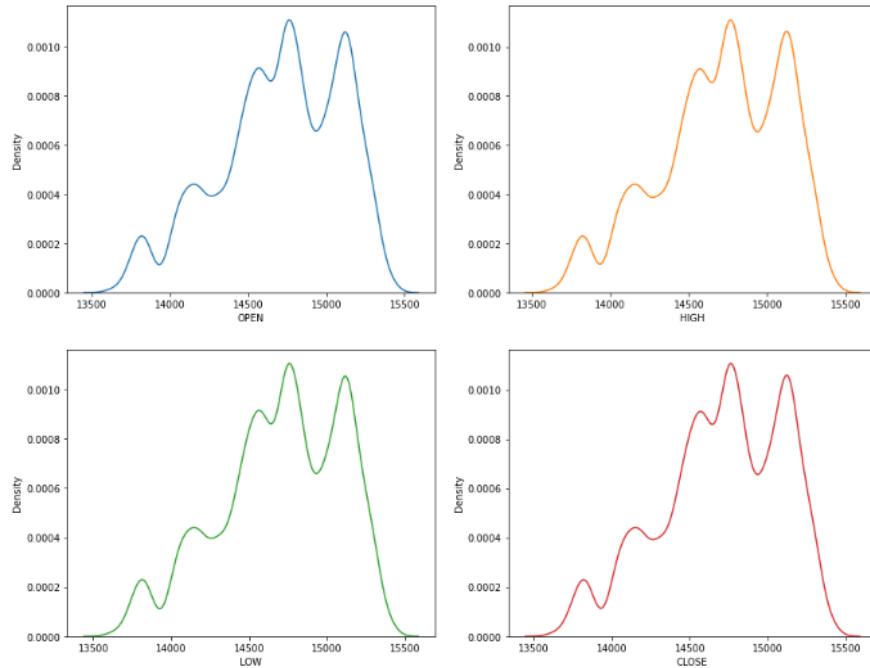
#### 1 d. Analyzing data

Data preparation is followed by data analysis, the discovery of hidden trends, and eventually the use of various machine learning models. A few features can be derived from the current features utilizing statistical techniques. Exploratory data analysis consists of various approaches, that is helpful for below points:

- Improve understanding of a data set,
- find underlying structure,
- extraction of key variables,
- spot outliers and anomalies,
- test underlying hypotheses,
- build parsimonious models, and

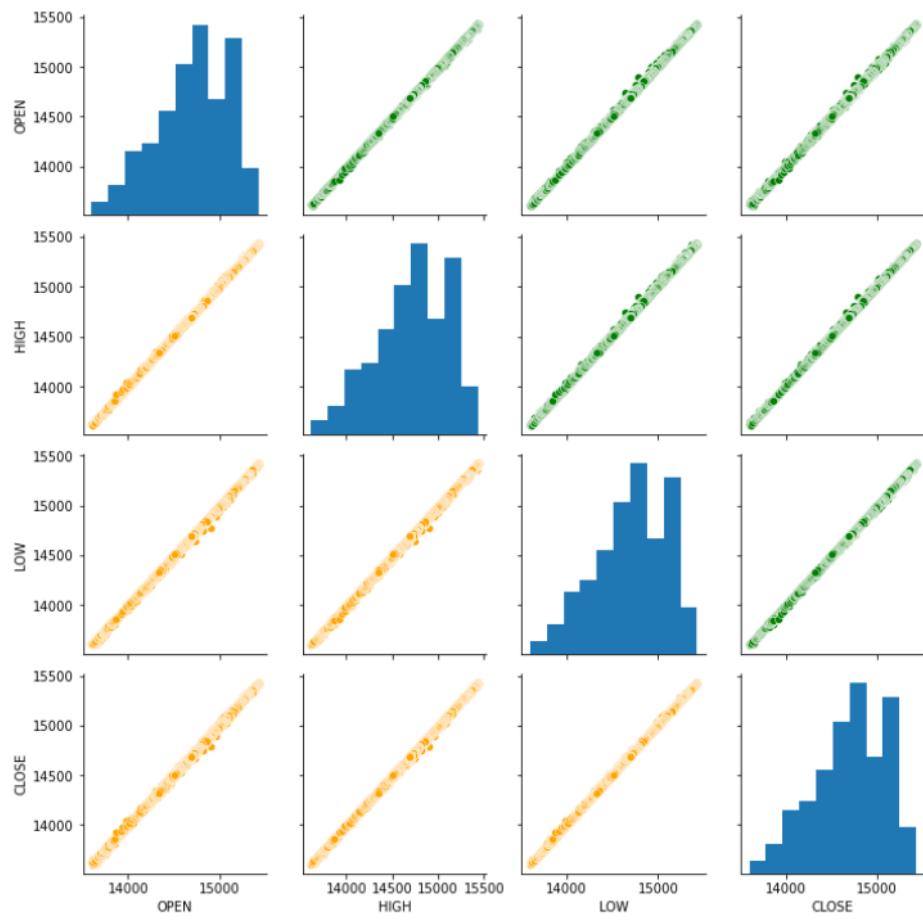
- choose the best factor settings.

Making a plot of the raw data, such as a block plot, probability plot, histogram, or bi-histogram. creating simple statistical plots from the raw data, such as mean plots, standard deviation plots, box plots, and major effects plots. putting numerous plots on a page to enhance our capacity for pattern identification when placing such plots.



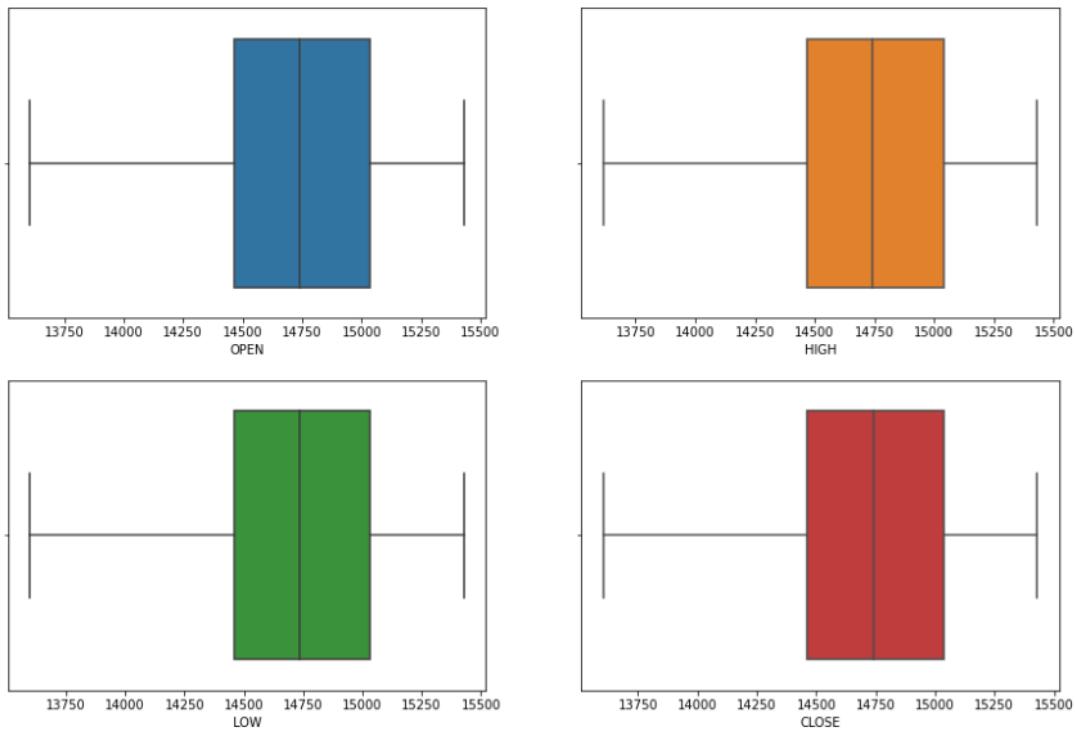
Distribution plots for the features

With the help of distribution plots, it has been identified that the data is skewed or normally distributed with different parameters.



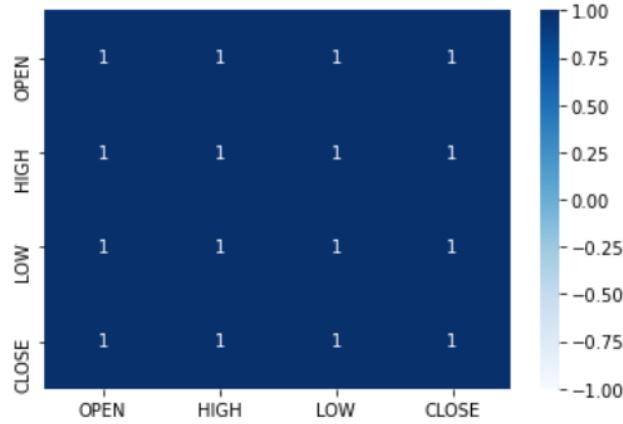
Scatter plots for the features

To identify the data distribution, scatter plots have been plotted.



Box plot to determine the outliers

Box plots are extensively used to determine the outliers in the dataset.



Heat map

The dataset's heat map/matrix examined the relationships between several characteristics. Now, the model has been trained using Linear, Poly SVM technique with 0.8 and 0.2 ratio with segregation of dependent and independent variables. The data has been transformed for feature scaling.

## Chapter 5: Machine Learning Model Performance

To predict the price of stock, a variety of machine learning methods have been available, including Support Vector Machine (SVM) as Linear, Support Vector Machine (SVM) as Poly, ANN, LSTM, Naive Bayes and Random Forest Algorithm. These models may be applied using the Python module Scikit Learn. R2, MAE, and MSE are a few of the parameters that are used to evaluate how well these models work. The formulae for these three parameters are as follows:

$$R^2 = 1 - \frac{\sum_{n=1}^{t=1} (y_i - \hat{y}_i)^2}{\sum_{n=1}^{t=1} (y_i - \bar{y})^2} \quad (1)$$

$$MAE = \frac{1}{n} \sum_{n=1}^{t=1} |y_i - \hat{y}_i| \quad (2)$$

$$MSE = \frac{1}{n} \sum_{n=1}^{t=1} (y_i - \hat{y}_i)^2 \quad (3)$$

Below are the models compared and these parameters calculated to get the best model

### 5.1- Random Forest Regressor

In order to produce more precise forecasts, this algorithm combines the less predictive algorithms. By merging the basic model, it creates a massive model. The features are sampled and transferred to the trees without replacement to produce highly uncorrelated decision trees. To select the best split, there must be less link between the trees. The crucial concept is aggregate uncorrelated trees, which set the random forest apart from the decision tree.

```
[ ] mse = mean_squared_error(rescaled_ytest,test_predict)
print('MSE: '+str(mse))
rmse = math.sqrt(mean_squared_error(rescaled_ytest,test_predict))
print('RMSE: '+str(rmse))
mape = np.mean(np.abs(test_predict - rescaled_ytest)/np.abs(rescaled_ytest))
print('MAPE: '+str(mape))

MSE: 565.2411820504213
RMSE: 23.77480140927409
MAPE: 0.015401272104545958
```

## 5.2 - SVM Regressor

Regression analysis and classification both employ the SVM supervised machine learning method. Since processing takes a long time, modest datasets are often used. It discovers the Hyperplane that divides the feature into many pieces. It provides an ideal hyperplane that categorizes many domains. Support vector points are the data points that are closest to the hyperplane, and margins are the separations between these points and the vector plane.

$$y = w_0 + \sum_{i=0}^m w_i x_i$$

The proposed work has exploited SVM for regression analysis. The performance depends on kernel function selection as a Non-parametric technique. Linear, Radial Basis Function and Polynomial are the kernels of support vector machine algorithms.

Support vector machines provide the following benefits:-

Effective in high-dimensional spaces.

SVMs are excellent when we don't know anything about the data. even with unstructured and semi-structured material, it performs effectively.

Even if there are more dimensions than samples, the method is still efficient.

It scales to high dimensional data rather well.

-It is memory efficient because the decision function uses a collection of training points known as support vectors.

Being adaptable allows for the specification of several kernel functions for the decision function. There are common kernels available, but you may also define your own kernels.

SVR Linear:

```
[ ] mse = mean_squared_error(y_test_e,y_test_pred_e)
print("Mean_squared_error is: ", mse)
r2score = r2_score(y_test_e , y_test_pred_e)
print("R2 score is: " , r2score)

Mean_squared_error is:  7971.497158328377
R2 score is:  0.9610811950203848
```

SVR Poly:

```
[ ] from sklearn.metrics import mean_squared_error, r2_score
mse = mean_squared_error(y_test_e,y_test_pred_f)
print("Mean_squared_error is: ", mse)
r2score = r2_score(y_test_e , y_test_pred_f)
print("R2 score is: " , r2score)

Mean_squared_error is:  15.522993756385382
R2 score is:  0.9999242129358256
```

### 5.3 Xgb Regressor

This approach makes use of gradient boosting to discover the function that fits the input data the best. By learning from the mistakes of prior decision tree training, it trains several decision trees and gets better each time. The goal of each iteration is to reduce the discrepancy between the anticipated and input data. When it comes to performance and scalability in a distributed context, it outperforms gradient boosting. Although accurate, it once again falls into the overfitting issue area, which may be resolved by adjusting the hyperparameter.

```
[ ] import xgboost
import lightgbm

gbm = lightgbm.LGBMRegressor()
gbm.fit(X_train_c, y_train_c)
prediction_c=gbm.predict(X_test_c)
print(mean_squared_error(prediction_c, y_test_c))

42.215051073569164
```

```
▶ xgb = xgboost.XGBRegressor()
xgb.fit(X_train_c, y_train_c)
prediction_c = xgb.predict(X_test_c)
print(mean_squared_error(prediction_c, y_test_c))
```

```
□ 51.6639762495133
```

## 5.4 - Experimental results with compare to auto h2o library

There is an open source library for python, that runs different algorithms and tries the best value for the hyperparameter to obtain the best results and then produce the chart based on it. It is based on the JVM and works as a distributed system with in memory machine learning techniques.

MSE: 52.61985338060575							
RMSE: 7.253954327165684							
MAE: 4.410230768008925							
RMSLE: 0.005876066245025321							
Mean Residual Deviance: 52.61985338060575							
R^2: 0.9997441716107782							
Null degrees of freedom: 986							
Residual degrees of freedom: 981							
Null deviance: 283702579.88870408							
Residual deviance: 51935.79528665787							
AIC: 6726.557938412846							
Cross-Validation Metrics Summary:							
mae	mean	sd	cv_1_valid	cv_2_valid	cv_3_valid	cv_4_valid	cv_5_valid
mean_residual_deviance	4.4064765	0.3122633	4.831698	4.0752583	4.627423	4.2255373	4.272467
rmse	52.98148	28.916622	103.4491350	32.865864	49.900578	41.76592	36.925907
null_deviance	40740516.0000000	5632848.0	40887808.0000000	32194294.0000000	39371752.0000000	44161680.0000000	47087044.0000000
r2	0.9997433	0.0001280	0.9995249	0.9998294	0.9997341	0.9998010	0.999827
residual_deviance	10387.159	5315.7344	19344.988	5488.5996	10429.221	8770.843	7902.1436
rmse	7.1014457	1.7856885	10.170995	5.732876	7.064034	6.4626555	6.0766687
rmsle	0.0056489	0.0018366	0.0085043	0.0038378	0.0059166	0.0057434	0.0042424

## Chapter 6 : Integrating with Project

The final step for the project is to integrate the model with the python project to give an interactive interface to the trader, so that it can abstract the complexity of the model and operate efficiently. After comparing the different models, the poly models proved to be best among the others. To save the model in a file, there is an external library available in python named “pickle”. With the help of pickle, the model can be saved and load from the binary file.

```
[ ] filename = 'svr_linear.pkl'  
pickle.dump(svr, open(filename, 'wb'))  
  
filename = 'svr_poly.pkl'  
pickle.dump(model_poly, open(filename, 'wb'))
```

Fig: Saving the model

Once the model has been saved to the file, it will be loaded into the python application and few of the parameters are taken as input from the user, like the name of stock where the user wants to trade with a limited of companies available, as the model needs to be trained for that company share price from the past. The second essential parameter is the amount of the money to invest in form of integer value, platform fees also play a vital role, and risk level basically decides if the user is willing to take the risk or not. Once the mandatory inputs are provided by the user, the program processes it using the model and produces the output that will help the trader to trade wisely for the stock.

```
amount = input('Enter money to invest in Rupees: ')
loss_margin = input('Enter how much loss can be acceptable in rupees?: ')
risk_level = input('Enter risk level (1 would be minimum and 5 would be maximum): ')
platform_fees = input('Enter the platform fees?: ')

## Fetch the last evening news to calculate sentiment score
news_data = read_news()
print('News data: ', news_data[0], news_data[1][0:100])

## Convert the news headlines into sentiment score using vader
sentiment_data_numerical = calculate_sentiment_data(news_data)

## load the model
stock_prediction_model = pickle.load(open('./models/stock_prediction.model', 'rb'))

## Predict the close value for the model by feeding all the inputs
```

main.py

C:\Users\ss\PycharmProjects\news\_for\_stocks\_dissertation\venv\Scripts\python.exe C:/Users/ss/PycharmProjects/news\_for\_stocks\_dissertation/main.py

Enter stock name: **HdfcBank**

Enter money to invest in Rupees: **5100**

Enter how much loss can be acceptable in rupees?: **100**

Enter risk level (1 would be minimum and 5 would be maximum): **3**

Enter the platform fees?: **100**

Reading past news

News data: 12-03-2019 US Federal Reserve keeps rates steady, signals no change in 2020, Assets of companies under insolven

Buy 4 stocks of HdfC

## **Chapter 7 : Conclusion and Inference**

This dissertation gives an understanding of implementation of machine learning techniques to identify the factors affecting the price of the stock and providing a confidence to the trader to invest in the stock market in the form of either buying or selling. The proposed project gives equal weightage to the trader capabilities of bearing the loss. The main objective of the project is to save the trader from having a loss and gaining a maximum profit on the good day. The proposed model is trained using the numerical parameters like the open, low, high as well as the textual parameters like the news headlines that is indirectly transformed into the numerical values using the sentiment analyser. The study has been done with various dataset like one year dataset and two year dataset and the final model is trained using the five year dataset. In order to provide better results, the accuracy and performance of various models are compared. Many prediction models are examined to see which one predicts the closing price of the stock most accurately. The model is then further integrated with the python project that gives an interactive interface to the traders to conveniently use it and get the results.

### **Directions for future work:**

Deep learning techniques can be implemented for improved results and it further helps to reduce one extra step for conversion of textual data into numerical data.

An enhanced user interface can be designed using web technologies that will help the end user to browse and select the options easily.

Proper way of handling the weekend data, as the stock market doesn't operate on weekends however the collective news is of three days, that may vary the predictions.

Grouping of potential stocks instead of calculating in isolated environments of the value of individual stocks.

## **Chapter 8 : References**

- [1] Thomas, J. D., and Sycara, K., 2000. Integrating Genetic algorithms and text learning for financial prediction. GECCO, July 8-12, Las Vegas, USA, pp. 72-75.
- [2] <https://economictimes.indiatimes.com/defaultinterstitial.cms>
- [3] <https://nsetools.readthedocs.io/en/latest/>
- [4] <https://www.sciencedirect.com/science/article/abs/pii/S0957417414004473>
- [5] NLP in Stock Market Prediction: A Review by Rodrigue Andrawos
- [6] Sheikh Abdullah, Mohammad Rahaman, and Mohammad Rahman. Analysis of stock market using text mining and natural language processing, 05 2013.
- [7] Jigar Patel et al. "Predicting stock and stock price index movement using trend deterministic data preparation and machine learning techniques". In: Expert systems with applications 42.1 (2015), pp. 259–268.
- [8] David MQ Nelson, Adriano CM Pereira, and Renato A de Oliveira. "Stock market's price movement prediction with LSTM neural networks". In: 2017 International joint conference on neural networks (IJCNN).
- [9] Stock Market Prediction Using LSTM Recurrent Neural Network  
<https://www.sciencedirect.com/science/article/pii/S1877050920304865>
- [10] Vader Sentiment Analyser C.J. HuttoEric Gilbert  
<https://ojs.aaai.org/index.php/ICWSM/article/view/14550>
- [11] [https://www.researchgate.net/publication/360726066\\_NLP\\_in\\_Stock\\_Market\\_Prediction\\_A\\_Review](https://www.researchgate.net/publication/360726066_NLP_in_Stock_Market_Prediction_A_Review)  
NLP in Stock Market Prediction: A Review by Rodrigue Andrawos

## Appendix

The GitHub code utilizes the trained model created using collab notebook

**GitHub:**

<https://github.com/sarvsav/dissertation>

**Collab:**

<https://colab.research.google.com/drive/1DXR8oXGeJ7yX96o72R2uoPWJGffi-hDj?usp=sharing>

**Dataset:** Uploaded in github repository inside data folder

(<https://www.kaggle.com/datasets/debashis74017/nifty-50-minute-data>)

**Models:** Uploaded in github repository inside models folder

## 1 Checklist of items for the Final Dissertation Report

a)	Is the Cover page in proper format?	Y / N
b)	Is the Title page in proper format	Y / N
c)	Is the Certificate from the Supervisor in proper format?  Has it been signed?	Y / N
d)	Is Abstract included in the Report?  Is it properly written?	Y / N
e)	Does the Table of Contents page include chapter page numbers?  1	Y / N
f)	Does the Report contain a summary of the literature survey? Are  i) the Pages numbered properly?  ii) Are the Figures numbered properly?  iii) Are the Tables numbered properly?  ) Are the Captions for the Figures and Tables proper?  iv) Are the Appendices numbered?  v)	Y / N Y / N Y / N Y / N Y / N
g)	Does the Report have Conclusion / Recommendations of the work?	Y / N
h)	Are References/Bibliography given in the Report?	Y / N
i)	Have the References been cited in the Report?	Y / N
j)	Is the citation of References / Bibliography in proper format?	Y / N

### Declaration by Student:

I certify that I have properly verified all the items in this checklist and ensure that the reports are in proper format as specified in the course handout.

Place: Gurgaon

Signature of the student

Name: Sarvsav Sharma

ID No: 2020SC04239

## ORIGINALITY REPORT



## PRIMARY SOURCES

- |   |   |     |
|---|---|-----|
| 1 | Submitted to Birla Institute of Technology and Science Pilani | 8%  |
|   | Student Paper   |     |
| 2 | arxiv.org   | 1 % |
|   | Internet Source   |     |
| 3 | kth.diva-portal.org   | 1 % |
|   | Internet Source   |     |
| 4 | Submitted to Universität Liechtenstein                        | 1 % |
|   | Student Paper   |     |
| 5 | Submitted to Liverpool John Moores University                 | 1 % |
|   | Student Paper   |     |
| 6 | papasearch.net  | 1 % |
|   | Internet Source   |     |
- 

Exclude quotes      On

Exclude bibliography      On

Exclude matches      < 1%