

“Stock Market Prediction Using ML”

Project report submitted in partial fulfillment of the requirement for the
degree of Bachelor of Technology

in

Computer Science and Engineering/Information Technology

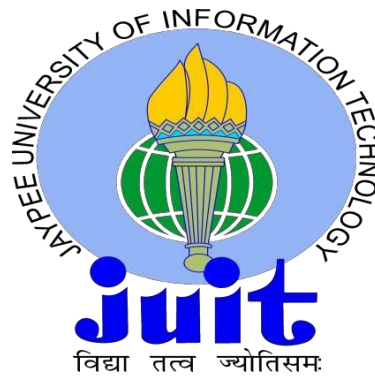
By

Dhairya Agarwal(161462)

Under the supervision of

Dr. Ruchi Verma

to



Department of Computer Science & Engineering and Information
Technology

**Jaypee University of Information Technology Wahnaghat, Solan-
173234, Himachal Pradesh**

Candidate's Declaration

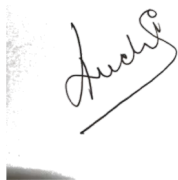
We hereby claim that the work that we are representing in this report entitled “Framework on Automated Trade Systems using Time-Series Data and ML Classifiers” in partial fulfillment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science and Engineering** submitted in the department of Computer Science & Engineering and Information Technology, Jaypee University of Information Technology Waknaghat is an authentic record of my own work carried out over a period from July 2019 to December 2020 under the supervision of Dr. Ruchi Verma Assistant Professor(Senior Grade).

The matter embodied in the report has not been submitted for the award of any other degree or diploma.

Dhairya Agarwal

Dhairya Agarwal, 161462

This is to certify you that the above statement made by the students are true to the best of my knowledge.



Dr. Ruchi Verma

Assistant Professor(Senior Grade)

Computer Science Department

Dated:

Acknowledgement

I would like to express my profound appreciation to all those who provided us the possibility to complete this report. A special gratitude we give to our final year project supervisor, Dr. Ruchi Verma, whose contribution in stimulating suggestions and encouragement, helped me and my partner to coordinate our project well especially in writing this report.

Furthermore we would also like to acknowledge with much appreciation the crucial role of Jaypee University of Information Technology, who gave the permission to use all the required equipment and the necessary materials to complete the project And framework on automated trade system using time-series data and ML classifiers .A special thanks goes to my supervisor, Dr. Ruchi Verma, who help me to assemble the parts and gave suggestion about the project “Algorithmic Trading” he have been invested his full effort in guiding us for achieving the goal. We have been to appreciate the guidance given by other supervisor as well as the panels especially in our project presentation that has improved our presentation skills thanks to their comment and advices.

Thanking you,

Dhairya Agarwal(161462)

Abstract

Determining the stock price in the market on any given day is a very complicated task even for the experts on the field as it depends on the large number of factors then human mind can process and that is where ML comes in the recent years ML has grown to a point where it can take in a large amount of data and find out patterns in that data which makes it very useful for stock price prediction.

Generally, budgetary exchange gauge is an uncommonly jumbling framework, to control stocks as showed by your necessities, fuses awkward information of stocks and how these stocks can change their headways and by the entirety they will climb or down considering some money related conditions. Issue is that, can a Machine predict these progressions and devise a kind of exchanging technique as indicated by the given information utilizing specific AI models.

Various shippers would take a gander at various degrees of a specific methodology, two or three venders likely won't put trust in exchanging utilizing machines, as they trust display renders itself once there is an immense expansion or lows. Merchants may utilize various classes with various parameters, and endeavor to devise a system that best fits the given illuminating social event.

One approach to manage build up a technique is accepting that, in a market a few affiliations are fundamentally related, negatively related, and some probably won't be related in any capacity whatsoever. Utilizing the affiliations machine can choose a decision on what fundamentals the information has been given, given the Classifiers, in the wake of expelling highlight sets and mapping them to names, the Classifiers take those capacities and fit the given information.

Table of Contents

Title Page

Declaration of the Student

Certificate of the Guide

Acknowledgement

Abstract

Introduction

List of Figures

List of Graphs

1. Introduction

Problem Definition

Project Overview and Specification

Hardware Specifications

Software Specifications

2. Literature Survey

Existing Systems

Proposed System

Feasibility Study

3. System Analysis & Design

Requirement Analysis

Data Extraction

Data Manipulation

Preprocessing data for ML

What are Classifiers?

Different Types Of Classifiers

Performing ML

4. Results and Outputs

5. Conclusions and recommendations

References

Appendices

List of Figures

1. Figure 1 - Code to Grab S&P 500tickers
2. Figure 2 - Output Of Grabbed 500tickers
3. Figure 3 - Code to grab stock data from Morningstar.
4. Figure 4 - Output Stock data of companies.
5. Figure 5 - Code to compile all close index of company in one data frame.
6. Figure 6 - Output close index of all companies together in one data frame.
7. Figure 7 - Code to find and visualize correlations.
8. Figure 8 - Output Of the correlation table.
9. Figure 9 - Heatmap of the correlations.
10. Figure 10 - Code to set trading conditions and data processing for labels.
11. Figure 11 - Code to extract feature sets and map them to labels.
12. Figure 12 - Pie chart
13. Figure 13 - Layer diagram
14. Figure 14 - ReLU function
15. Figure 15 - Neural network
16. Figure 16 – Recurrenet neural network
17. Figure 17 - Code of implementing Classifiers and performing ML.
18. Figure 18 - Output data spread and predicted spread
19. Figure 19 - Graph Of MMM company that year.

Chapter 1

Introduction

Problem Definition:

Stock market checks are an unprecedented spelling work, embodied in the high substance of budgetary exchange limit, and thinking about different circumstances can provoke market dissatisfaction. While some dealer may fight that the market and is fair in itself, and if there is a new check or someone in the market that collects it from the standard, they themselves review it and charge it, similarly.

Think of the masters of wealth, think of animation, buy low, move high, yet don't provide enough setting to make decisions about fitting efforts. Before an inspector holds assets in any stock, he must identify how the cash market continues. Keeping assets in an amazing stock in spite of terrible times can have been terrible consequences, while in preferred times the importance for a common stock may be under the core interests of the central government. . Today's cash-related Money Stars are moving towards this issue of trading because they do not and the big bag is to buy related shares or which clearer than predicted on general explanations .

This application looks for the requirements and accommodation in classroom.

Project Overview and Specifications

Man-made thinking (AI) to accept a basic activity in our regular daily existence money related applications like trading, it is an improvement towards some other time of development. This endeavor includes an utilization of Artificial intelligence on budgetary algorithmic trading.

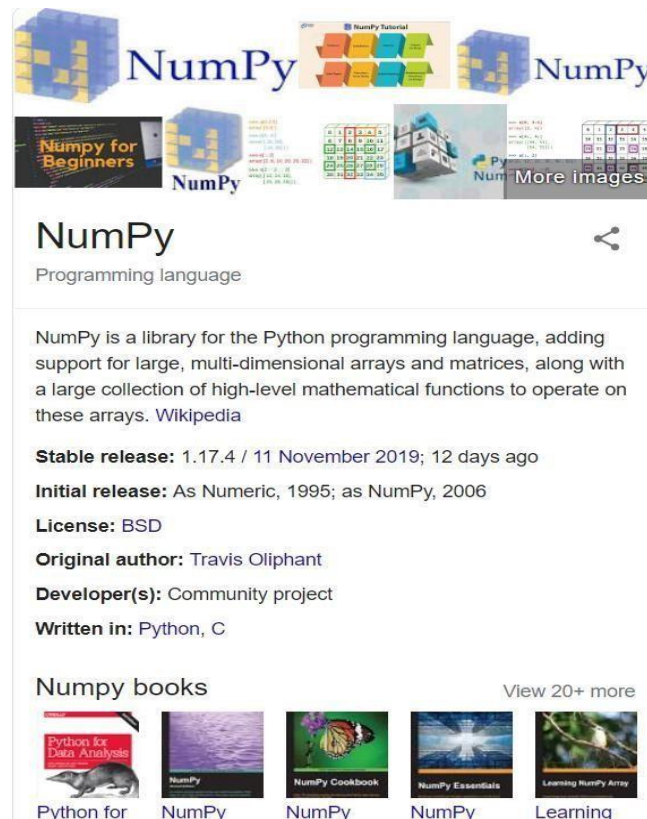
Motorized trading structures incorporates to choose extremely speedy .

Artificial intelligence is a subset of AI and all around gives game plans which gain for a reality without being unequivocally altered.

Hardware Specifications

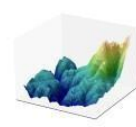
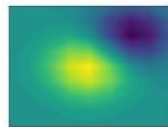
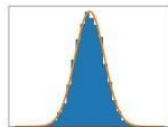
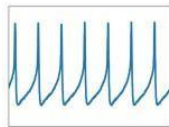
Automation of models may sound too much incredible yet requires relatively better than average PC with a not too bad editor and you're good to go, next to no need of extra gear judgments.

Software Specifications



The image shows a screenshot of a project page for NumPy. At the top, there are several images related to NumPy, including the logo, a tutorial diagram, and a book cover titled 'NumPy for Beginners'. Below the images, the text reads 'NumPy' followed by 'Programming language'. A description states: 'NumPy is a library for the Python programming language, adding support for large, multi-dimensional arrays and matrices, along with a large collection of high-level mathematical functions to operate on these arrays. Wikipedia'. Below the description, there is a 'Stable release' section with the version '1.17.4 / 11 November 2019; 12 days ago', an 'Initial release' section with 'As Numeric, 1995; as NumPy, 2006', a 'License' section with 'BSD', an 'Original author' section with 'Travis Oliphant', a 'Developer(s)' section with 'Community project', and a 'Written in' section with 'Python, C'. At the bottom, there is a section titled 'Numpy books' with a 'View 20+ more' link. Below this section, there are five book covers: 'Python for Data Analysis', 'NumPy', 'NumPy Cookbook', 'NumPy Essentials', and 'Learning NumPy Array'.

Matplotlib is a Python 2D plotting library which produces publication quality figures in a variety of hardcopy formats and interactive environments across platforms. Matplotlib can be used in Python scripts, the Python and IPython shells, the Jupyter notebook, web application servers, and four graphical user interface toolkits.



Matplotlib tries to make easy things easy and hard things possible. You can generate plots, histograms, power spectra, bar charts, errorcharts, scatterplots, etc., with just a few lines of code. For examples, see the [sample plots](#) and [thumbnail gallery](#).

For simple plotting the `pyplot` module provides a MATLAB-like interface, particularly when combined with IPython. For the power user, you have full control of line styles, font properties, axes properties, etc, via an object oriented interface or via a set of functions familiar to MATLAB users.

PANDAS

```
import pandas as pd
```

```
df = pd.read_csv("f500.csv", usecols = ["company", "rank", "revenues"])
```

```
df
```

	company	rank	revenues
0	Walmart	1	485873
1	State Grid	2	315199
2	Sinopec Group	3	267518

```
import pandas as pd
import datetime
import pandas_datareader.data as web
from pandas import Series, DataFrame
```

Chapter 2

Literature Survey

Existing Systems

Cash related business choices have been a solid and di ment driven task budgetary information. Figure about exchanging of affirmations high exactness amendment returns gains for stock agents. In building up in the hold of budgetary business money related information, the headway of the benefit model for estimation closes is very faction, and one ought to be cautious. The thought attempted to make a model for the protection business theory and to decide to purchase/hold stocks utilizing information mining and AI systems. Computer based intelligence structures, for example, k-closest neighbors (k-NN).

The timetable sets on the information and using it like the theory model. Long lasting guaranteed exchanging information have been been utilized for stock sign checks. In context on edifying characterization models may demonstrate purchase/hold for money related exchange the type of yields. The standard this task is to pass on the yield signal (purchase/h0ld) as showed by clients, for example, central commitment, time of exertion, least favored area, most clear trouble, information mining and utilization of AI structure do.

Envisioning systems for stock spending is the subject of a wide-going dialog in different fields, including exchanging, records, bits of information, and programming building. In the preservation exchange, benefactors can sell their produce by charge or sell their exertion in the event that they can pick when to enter and leave a position. Commonplace brokers regularly utilize the key just as an exceptional test for esteeming partakes in picking experience choices. The essential test incorporates an assessment of affiliation basics, for instance, salary and advantages, office status, improvement rates, and then some. Explicit appraisal depends, obviously, on the assessment of more seasoned respect developments. Because of the exhibiting powers, the economy will look for progression over wide stretches ever and one instance of suffocation. Stocks exchange a typical circumstance where the economy works from a period.

Stood separated from the present work, this exertion endeavors stock exchanging choices utilizing the regular direct of exchanging structures inside the foundation of money related fiscal and business positions.

The target work is to improve the medium to long haul for S&P 500 cash related professional resources. The wells of the information are unique sign information and budgetary marker information. The three models are then used to anticipate purchasing/selling choices..

Proposed Framework

As the exchange slices identified with overhead cash are discussed, this is an exceptionally enormous point and a piece everything being equal and the craving is to be certain how well the model partner fits for a given dataset May and whether it seriously portrayed the outcomes and assesses accomplish. No. However each model gives some impact to all the concentration and reason, every one of them requiring a silly blueprint of the relationship of any stock trade to look at game-play parts paying little respect to which one Can.

Request to set up datasets for extra AI cohorts who will in the long run theory the etching and pass the immeasurably significant great approach to discover the likelihood of time) and control the information in the model on the three essential places of the market to purchase, hold, Were. Also, the security for each relationship against their tickers for selling and doing so was overseen in the business information machine

Probablity Study

The common sense of the overhead model will be plaid enabled to plaid and be fabricated near the diagram of positive connections for that timeframe and see the model.

As a future degree in our undertaking, we will likewise utilize Quantopiana to assemble ways to deal with the web based exchanging stage and test them back.

Chapter 3

System Analysis and Design

Technical Analysis

Specific detailing is useful for measuring future monetary stock growth dependent on stock authentic growth. Specific sanctions do not estimate stock costs, yet relying on outdated investigations, particular points of imprisonment can trace the development of stocks to the current economic situation after some time. Specific valuation helps the investor to estimate stock price growth (up / down) in the interim at that particular time. The generality of specialized assessment is a variety of modalities that indicate costs over a period.

The organization rundown from Wikipedia are being spared and the stock information is being intentionally eliminated in the logical inconsistency of each organization ticker.

At that point each organization's close file is considered and put into an information casing and an effort is made to discover a consistency between each organization and later on to prepare the information and the stock price, mass And depends on making specific special criteria based on the nearest value and the dependent. The development of costs will advance special meters that will help determine the objective rate of purchase, sale, hold.

Data Extraction

	Date	Open	High	Low	Last	Close	Total Trade Quantity	Turnover (Lacs)
0	2018-10-08	208.00	222.25	206.85	216.00	215.15	4642146.0	10062.83
1	2018-10-05	217.00	218.60	205.90	210.25	209.20	3519515.0	7407.06
2	2018-10-04	223.50	227.80	216.15	217.25	218.20	1728786.0	3815.79
3	2018-10-03	230.00	237.50	225.75	226.45	227.60	1708590.0	3960.27
4	2018-10-01	234.55	234.60	221.05	230.30	230.90	1534749.0	3486.05

To begin, organizations have been a basic basis for which we can gauge measurable ways for the requirement. Each organization will have been a personal record of stock information from 1/1/2000 to 31/12/2017. Right off the bat, a rundown of organizations was required, Random is mined from Wikipedia, there the S&P list is in a table configuration, whatever it is but hard to deal with.

Use random to turn the pickle off and if random changes, check it for a clear timeframe Change for Redeem the ticker random, so as not to hit Wikipedia more than riding content every time.

There are tickers of 520 organizations, a stock evaluating the information of each organization is required. The stock evaluating the information of the initial 18 organizations is theoretical, with each organization having stock information for each organization for approximately 6020 confirmations. Organizations that were started after 2020 and have been empty properties have been their gateway pressed by zero.

Select the information that panda-information peruser uses, a python mining library.

Select information from Morningstar and locally spares the information in the .csv design and the information will be useful for later effects.

Currently, stocks information , can move in the direction of information and realize which files are in our information.

```

gettingsnp500.py - C:\Users\Root\Documents\Project\gettingsnp500.py (3.6.5)
File Edit Format Run Options Window Help

import sys as bs
import datetime as dt
import matplotlib.pyplot as plt
from matplotlib import style
import numpy as np
import os
import pandas as pd
import pandas_datareader.data as web
import pickle
import requests

style.use('ggplot')

def save_sp500_tickers():
    resp = requests.get('https://en.wikipedia.org/wiki/List_of_S%26P_500_companies')
    soup = bs.BeautifulSoup(resp.text, "lxml")
    table = soup.find('table',{'class':'wikitable sortable'})
    tickers = []
    for row in table.findAll('tr')[1:]:
        ticker = row.findAll('td')[0].text
        tickers.append(ticker)

    with open("sp500tickers.pickle","wb") as f:
        pickle.dump(tickers, f)

    print(tickers)

    return(tickers)

save_sp500_tickers()

```

Fig : Code to Grab S&P tickers

```

Python 3.6.5 Shell
File Edit Shell Debug Options Window Help
Python 3.6.5 (vs.6.5:f59c0932b4, Mar 29 2018, 16:07:46) [MSC v.1900 32 bit (Intel)] on win32
Type "copyright", "credits" or "license()" for more information.
>>>
===== RESTART: C:\Users\Root\Documents\Project\gettingsnp500.py =====
[['MMM', 'ABT', 'ABBV', 'ABMD', 'ABN', 'ATVI', 'ADBE', 'AMD', 'AAP', 'AES', 'AET', 'AMG', 'AFL', 'A', 'APD', 'AXAM', 'ALX', 'ALB', 'ARE', 'ALXN', 'ALGN', 'AL
LE', 'AGN', 'ADS', 'LNT', 'ALL', 'GOOGL', 'GOOG', 'MO', 'AMZN', 'AEP', 'DAL', 'AEP', 'AXP', 'AIG', 'AMT', 'AWK', 'AMP', 'ABC', 'AME', 'AMGN', 'ARH', 'APC',
'ADI', 'ANSS', 'ANTM', 'AON', 'AOS', 'AFA', 'AIV', 'AAPL', 'AMAT', 'AFTV', 'ADM', 'ARNC', 'ANET', 'AJG', 'AIZ', 'T', 'ADSK', 'ADP', 'AZO', 'AVB', 'AVY', 'BH
GE', 'BLI', 'BAC', 'BK', 'BAK', 'BBI', 'BDX', 'BRK-B', 'BBY', 'BIIB', 'BLK', 'HRB', 'BA', 'BKNG', 'BWA', 'BXP', 'BSX', 'BHE', 'BMY', 'AVGO', 'BR', 'BF-B', '
CHW', 'COG', 'CDNS', 'CEB', 'COF', 'CAH', 'CRX', 'CCL', 'CAT', 'CBSE', 'CBRE', 'CBS', 'CELG', 'CNC', 'CNP', 'CTL', 'CERN', 'CF', 'SCHW', 'CHTR', 'CVX', 'C
G', 'CB', 'CHD', 'CI', 'XEC', 'CINF', 'CTAS', 'CSCO', 'C', 'CFG', 'CTXS', 'CLX', 'CME', 'CMS', 'RO', 'CTSH', 'CL', 'CMCSA', 'CMA', 'CAG', 'CXO', 'COF', 'ED
G', 'STZ', 'COO', 'CPRT', 'GLW', 'COST', 'COTY', 'CCI', 'CSX', 'CHI', 'CVS', 'DHI', 'DHR', 'DRI', 'DVA', 'DE', 'DAL', 'XRAY', 'DWN', 'DLR', 'DFS', 'DISCA', 'D
ISCR', 'DISH', 'DG', 'DLR', 'D', 'DOV', 'DND', 'DTE', 'DRE', 'DUK', 'DXC', 'ETFC', 'EMN', 'ETN', 'EBAY', 'ECL', 'EIX', 'EW', 'EA', 'EHR', 'ETR', 'EOG', 'E
FX', 'EQIX', 'EQR', 'ESS', 'EL', 'EVRG', 'ES', 'RE', 'EXC', 'EXPE', 'EXPD', 'ESRX', 'EXR', 'XOM', 'FFIV', 'FB', 'FAST', 'FRT', 'FDX', 'FIS', 'FITB', 'FE', '
FISV', 'FLT', 'FLIR', 'FIS', 'FLR', 'FMC', 'FE', 'F', 'FTNT', 'FTV', 'FBHS', 'BEN', 'FCX', 'GFS', 'GRMN', 'IT', 'GD', 'GE', 'GIS', 'GM', 'GPC', 'GILD', 'GPN
', 'GS', 'GT', 'GWW', 'HAL', 'HBI', 'HOG', 'HRS', 'HIG', 'HAS', 'HCA', 'HCP', 'HP', 'HSIC', 'HSY', 'HES', 'HPE', 'HLT', 'HFC', 'HOLX', 'HD', 'HON', 'HRL', '
RST', 'HPO', 'HUM', 'HEAN', 'HII', 'IDXX', 'INFO', 'IITW', 'IHM', 'IR', 'INTC', 'ICE', 'IRM', 'INCY', 'IP', 'IPG', 'IFF', 'INTU', 'ISRG', 'IVZ', 'IRFP', 'IQ
V', 'IRM', 'IQVY', 'JEC', 'JBT', 'JEF', 'SJM', 'JNJ', 'JCI', 'JPM', 'JNNA', 'KRO', 'K', 'KEY', 'KEYS', 'RMB', 'KIM', 'KMI', 'KLAC', 'KSS', 'KHC', 'KR', 'LE
', 'LEL', 'LH', 'LRCX', 'LEG', 'LEN', 'LIV', 'LNC', 'LIN', 'LQ', 'LMT', 'L', 'LOW', 'LNB', 'MBS', 'MCO', 'M', 'MRO', 'MEO', 'MBA', 'MCC', 'MEM', 'MAS', 'MA
', 'MAT', 'MCO', 'MCD', 'MCR', 'MDT', 'MRK', 'MET', 'MTD', 'MGM', 'KORS', 'MCHP', 'MU', 'MST', 'MBA', 'MHR', 'TAB', 'MDLZ', 'MNST', 'MCO', 'MS', 'MOS', 'MS
I', 'MSCI', 'MELI', 'NDAQ', 'NOV', 'NTR', 'NTAP', 'NFLX', 'NWL', 'NEX', 'NEM', 'NWSA', 'NWS', 'NEE', 'NLSN', 'NKE', 'NI', 'NBL', 'JWN', 'NSC', 'NTRS', 'NOC
', 'NCLH', 'NRG', 'NUV', 'NUVA', 'ORLY', 'ONX', 'OMC', 'OKE', 'ORCL', 'PCAR', 'PKG', 'PH', 'PAIX', 'PYEL', 'PNR', 'PBT', 'PEP', 'EKI', 'PRGO', 'PFE', 'PG',
'PM', 'PSX', 'PWR', 'PNC', 'RL', 'PFS', 'PEP', 'PEG', 'PG', 'PGR', 'BLD', 'PRU', 'PEG', 'PSA', 'PHM', 'FVH', 'QROV', 'PWR', 'QCOM', 'DGX', 'RJT', 'RT
N', 'O', 'RHT', 'REG', 'REGN', 'RE', 'RSG', 'RMD', 'RHI', 'ROK', 'COL', 'ROL', 'ROP', 'ROST', 'RCL', 'CRM', 'SBAC', 'SOG', 'SLB', 'STX', 'SEE', 'SRE', 'SHW
', 'SPG', 'SWKS', 'SLG', 'SNA', 'SO', 'LUV', 'SPGI', 'SWK', 'SBOX', 'STI', 'SRCL', 'SYR', 'STI', 'SIVB', 'SYMC', 'SVT', 'SNPS', 'SY', 'TRON', 'ITWO', 'TRR
', 'TGT', 'TEL', 'FTI', 'TXN', 'ITX', 'TMO', 'TIF', 'TWTX', 'TXN', 'TME', 'TSS', 'TSCO', 'TDG', 'TRV', 'TRIP', 'FOXA', 'FOX', 'TSN', 'UDR', 'ULTA', 'USB', 'URR
', 'UA', 'UNE', 'UAL', 'UNH', 'UES', 'URI', 'UTX', 'URS', 'URN', 'VFC', 'VLO', 'VAR', 'VTR', 'VRSN', 'VRSK', 'VZ', 'VRTX', 'VIAB', 'V', 'VNO', 'VNC', 'WMT
', 'WBA', 'DIS', 'WM', 'WAT', 'WEC', 'WCG', 'WFC', 'WELL', 'WDC', 'WU', 'WRK', 'WY', 'WHR', 'WMB', 'WLTW', 'WINN', 'XEL', 'XRX', 'XLNX', 'XYL', 'YUM', 'ZBH', '
ZION', 'ZTS']

```

Fig: Output of Grabbed tickers

```

gettingsnp500.py - C:\Users\Root\Documents\Project\gettingsnp500.py (3.6.5)
File Edit Format Run Options Window Help

def get_data_from_morningstar(reload_sp500=False):
    if reload_sp500:
        tickers = save_sp500_tickers()
    else:
        with open("sp500tickers.pickle","rb") as f:
            tickers = pickle.load(f)

    if not os.path.exists('stock_dfs'):
        os.makedirs('stock_dfs')

    start = dt.datetime(2000,1,1)
    end = dt.datetime(2017,12,31)

    for ticker in tickers[:20]:
        print(ticker)
        str1 = "new"
        ticker_append = ticker + str1
        if not os.path.exists('stock_dfs/{}.csv'.format(ticker)):
            df = web.DataReader(ticker, 'morningstar', start, end)
            df.to_csv('stock_dfs/{}.csv'.format(ticker))
            df=pd.read_csv('stock_dfs/{}.csv'.format(ticker))
            df_col = ['Date', 'Open', 'High', 'Low', 'Volume', 'Close']
            new_df=df[df_col]
            new_df.to_csv('stock_dfs/{}.csv'.format(ticker_append), index=False)

        else:
            df=pd.read_csv('stock_dfs/{}.csv'.format(ticker))
            print(df.head())
            print('Already Have{}'.format(ticker))

get_data_from_morningstar()

```

Fig : Code to Grab stock data

```

Python 3.6.5 Shell
File Edit Shell Debug Options Window Help

MMM
Symbol      Date      Close    High     Low      Open     Volume
0  MMM  2000-01-03  47.188  48.250  47.032  48.03  2173400
1  MMM  2000-01-04  45.313  47.407  45.313  46.44  2713800
2  MMM  2000-01-05  46.625  48.125  45.563  45.57  3699400
3  MMM  2000-01-06  50.375  51.250  47.157  47.16  5975800
4  MMM  2000-01-07  51.375  51.907  49.969  50.57  4101200
Already HaveMMM
ABT
Symbol      Date      Close    High     Low      Open     Volume
0  ABT  2000-01-03  15.6639  16.1114  15.5520  15.7753  10667432
1  ABT  2000-01-04  15.3003  15.5520  15.1045  15.4121  9315154
2  ABT  2000-01-05  15.1886  15.6361  15.0208  15.0208  11758286
3  ABT  2000-01-06  15.4401  15.7758  15.1324  15.1324  13389878
4  ABT  2000-01-07  16.0555  16.2233  15.4401  15.4401  14315159
Already HaveABT
ABEV
Symbol      Date      Close    High     Low      Open     Volume
0  ABEV  2012-12-10  35.00  37.00  34.91  37.00  749378
1  ABEV  2012-12-11  35.00  35.00  35.00  35.00  0
2  ABEV  2012-12-12  33.36  35.25  33.09  35.25  2530442
3  ABEV  2012-12-13  33.80  34.24  32.61  33.50  4253303
4  ABEV  2012-12-14  33.00  34.08  32.65  33.90  4006673
Already HaveABEV
ABMD
Symbol      Date      Close    High     Low      Open     Volume
0  ABMD  2000-01-03  18.250  18.657  18.250  18.52  185600
1  ABMD  2000-01-04  17.813  18.500  17.000  18.50  34400
2  ABMD  2000-01-05  18.000  18.188  16.938  17.07  122800
3  ABMD  2000-01-06  18.032  18.063  17.625  17.63  84200
4  ABMD  2000-01-07  17.938  18.250  17.563  18.00  69000
Already HaveABMD
ACN
Symbol      Date      Close    High     Low      Open     Volume
0  ACN  2001-07-19  15.17  15.29  15.00  15.10  33703500
1  ACN  2001-07-20  15.01  15.05  14.80  15.05  9238500
2  ACN  2001-07-23  15.00  15.01  14.55  15.00  7501000
3  ACN  2001-07-24  14.86  14.97  14.70  14.95  3537300
4  ACN  2001-07-25  14.95  14.95  14.65  14.70  4208200
Already HaveACN
ATVI
Symbol      Date      Close    High     Low      Open     Volume
0  ATVI  2000-01-03  1.3699  1.3748  1.1664  1.3149  7226760

```

Fig : Output Stock data

Data Manipulation

At last that, stock considering information of affiliations is overseen, join this information in a solitary information plot. The aggregate of the stock records really beginning at now go with: Open, High, Low, Close, Volume.

In a brief range, check if any beguiling relationship data is found. To do this, imagine it's a huge proportion of data. Use library to endeavor this which can plot multiple graphs.

'gg plot' is an association is central, perform it by in a general sense including standard join _data diagram.

Direct, by making this craftsmanship, heatmap, this is an astoundingly clearly should change the veritable data.

Warmth map is gotten from the c-map, use RdY1Gnm, which is a covering - record that goes from red will give for expulsions attempts, green for positive affiliations. Included outlining bar that is a covering obliging. As time goes on, standard x and y center point shakes so there is a way to deal with oversee regulate supervise direct observe which concerns are which, since there is on an essential level plots.

To pass on express space between both the axes in 2d plane. Evaluations to make plots all around cautiously quiet to study, regardless, for this condition, it doesn't. By then also impudent the x center to be at the for all intents and purposes a relationship .

For this condition, the particular same structure works well.

Other than turn the x ticks, which are just the particular tickers, since incessantly they'll be turned out on a level plane. There are in excess of 500 checks turn far in the past a framework that will be irrationally epic to truly watch everything zoomed out.

There are a not a lot of that are unequivocally related there are some that are ridiculously blue down and some are not related at all by any stretch of the imagination. Looking affiliations, see that there are various affiliations. By a wide edge most by a long shot of affiliations are, expectedly plainly related.

```

gettingsnp500.py - C:\Users\Root\Documents\Project\gettingsnp500.py (3.6.5)
File Edit Format Run Options Window Help

def compile_data():
    with open("sp500tickers.pickle", "rb") as f:
        tickers = pickle.load(f)

    main_df = pd.DataFrame()
    count_ = 0
    for count, ticker in enumerate(tickers):
        count_ = count + 1;
        if count_ >= 20:
            break;
        else:
            str1 = "new"
            ticker_append = ticker + str1
            #print(ticker_append)
            df = pd.read_csv('stock_dfs/{}.csv'.format(ticker_append))
            df.set_index('Date', inplace=True)

            df.rename(columns = {'Close': ticker}, inplace=True)
            df.drop(['High', 'Low', 'Open', 'Volume'], 1, inplace=True)

            if main_df.empty:
                main_df = df
            else:
                main_df = main_df.join(df, how='outer')

        if count % 10 == 0:
            #print(count)

            #print(main_df.head())
            main_df.to_csv('sp500_joined_closes.csv')

compile_data()

```

Fig : Code to combine close indexes of company in one data frame.

```

3  ALXN  2000-01-06  7.7500  7.7500  7.2500  7.2500  800400
4  ALXN  2000-01-07  8.2500  8.3125  7.3750  7.5000  749200
Already HaveALXN
      Date      MMM    ABT    ABBV  ...    AMD    AAP    AES    AET
4690 2017-12-25  234.73  56.93  98.21  ...   10.54  100.55  10.71  179.96
4691 2017-12-26  235.45  57.00  97.75  ...   10.46  101.96  10.64  180.42
4692 2017-12-27  236.20  57.47  98.09  ...   10.53  99.77  10.67  180.85
4693 2017-12-28  235.72  57.46  97.79  ...   10.55  99.71  10.76  181.23
4694 2017-12-29  235.37  57.07  96.71  ...   10.28  99.69  10.83  180.39

```

Fig : Output in one data frame.

```

gettingsnp500.py - C:\Users\Root\Documents\Project\gettingsnp500.py (3.6.5)
File Edit Format Run Options Window Help
main_df = main_df.join(df_new, subset=
    if count % 10 == 0:
        #print(count)
        #print(main_df.head())
        main_df.to_csv('sp500_joined_closes.csv')
compile_data()

def visualize_data():
    df = pd.read_csv('sp500_joined_closes.csv')
    print(df.tail())
    df_corr=df.corr()
    print(df_corr.head())

    data = df_corr.values
    fig = plt.figure()
    ax = fig.add_subplot(1,1,1)

    heatmap = ax.pcolor(data , cmap=plt.cm.RdYlGn)
    fig.colorbar(heatmap)
    ax.set_xticks(np.arange(data.shape[0]) + 0.5, minor=False)
    ax.set_yticks(np.arange(data.shape[1]) + 0.5, minor=False)
    ax.invert_yaxis()
    ax.xaxis.tick_top()

    column_labels = df_corr.columns
    row_labels = df_corr.index

    ax.set_xticklabels(column_labels)
    ax.set_yticklabels(row_labels)
    plt.xticks(rotation=90)
    heatmap.set_clim(-1,1)
    plt.tight_layout()
    plt.show()

visualize_data()

```

Fig : Code to visualize correlations.

```

4690 2017-12-25 234.73 56.93 98.21 ... 10.54 100.55 10.71 179.96
4691 2017-12-26 235.45 57.00 97.75 ... 10.46 101.96 10.64 180.42
4692 2017-12-27 236.20 57.47 98.09 ... 10.53 99.77 10.67 180.85
4693 2017-12-28 235.72 57.46 97.79 ... 10.55 99.71 10.76 181.23
4694 2017-12-29 235.37 57.07 96.71 ... 10.28 99.69 10.83 180.39

[5 rows x 12 columns]
      MMM      ABT      ABBV      ...      AAP      AES      AET
MMM  1.000000  0.925710  0.921983  ...  0.859452 -0.304913  0.965608
ABT  0.925710  1.000000  0.914364  ...  0.879624 -0.217926  0.921008
ABBV 0.921983  0.914364  1.000000  ...  0.289796 -0.253843  0.907987
ABMD 0.866284  0.795963  0.834205  ...  0.715834 -0.074593  0.885329
ACN  0.953973  0.942296  0.839797  ...  0.898378 -0.043713  0.945445

[5 rows x 11 columns]

```

Fig 8: Output Of correlation table.

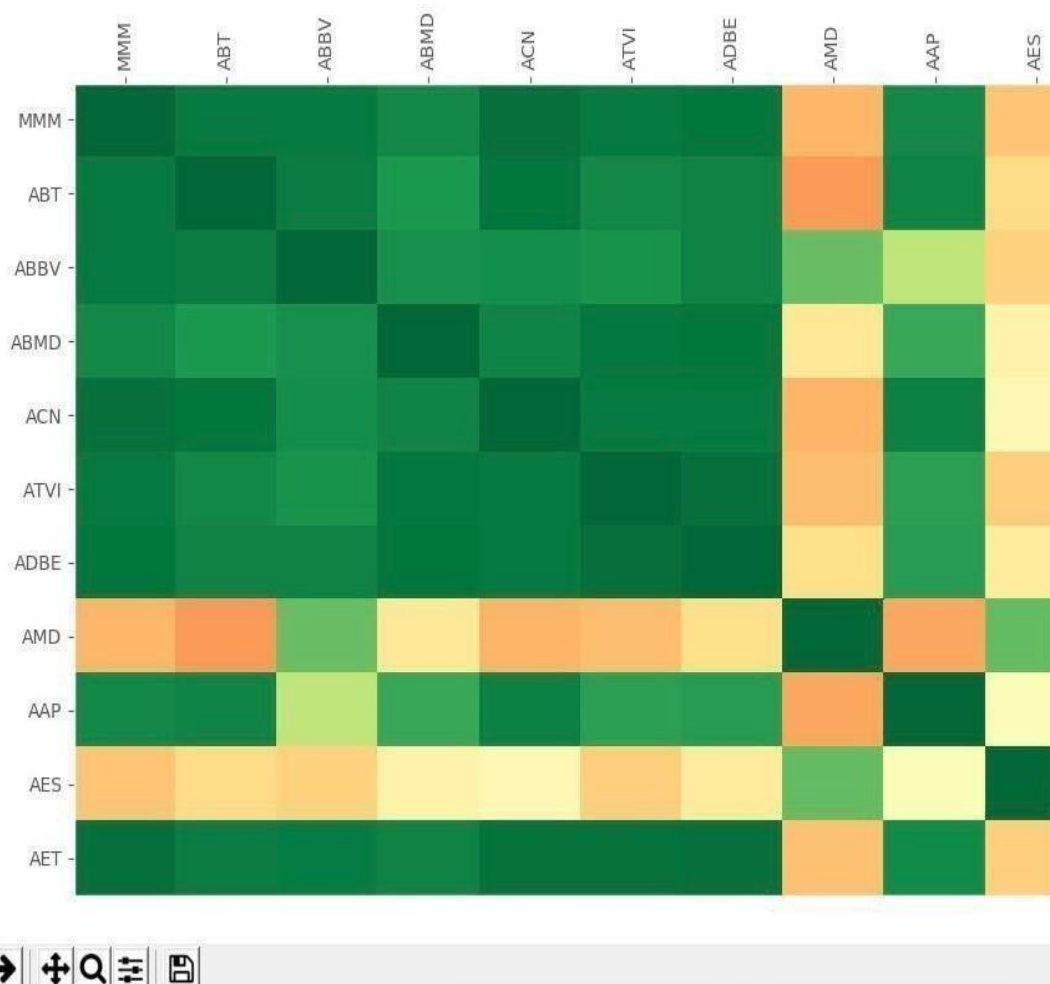


Fig : Heatmap of correlations.

Preprocessing the Data to anticipate ML

First thing this is one of the colossal strides of information mining what must be done before information mining so constantly information mining contains a monstrous proportion of joins which are utilized so as to mine information and get structures consequently this is a course of action of keen estimations which are being filtered through by scientist and for the most part when we consider information mining we basically think about these calculations in any case everything considered a gigantic measure of work models, etc yet routinely we imagine that this reason for get together of the work which is really everything looked at wires as a couple of stages a bit of the time information isn't on and on present in a solitary spot it may must be amassed from heaps of better places for instance on the off chance that you are doing some web application so the information may must be amassed from a couple of unquestionable region and accumulated condition in the event that you are the proprietor of recognize a chain of expect markets or a bit of that branch and what you need is to hide away this information into a particular zone so what you need is information get-together or information mix from various sources into a particular vault and that once in a while information so after you gather the information you need to clean the information directly at present clearing radiates an impression of being an unremarkable sort of less an amazingly manual concentrated work so after there is still some manual exertion yet there are do information collecting and cleaning you need to change over the information into a reasonable diagram and the union

change relies upon what sort of experience of mining should dismiss certain bits of the information certain properties you ought to consolidate just certain credits you should join certain qualities into different attributes and some time later utilize the ensuing plan has guarantee to the information mining figurings at last you model discover is a great deal of gigantic alright so you should deal with the entire of the qualities that happen for this trademark into the range address 0 to 1 so there will be fragmentary attributes some spot in the level of 0 and 1 alright so in all to do that plainly for this situation is pivotal you essentially need to oblige the aggregate of the open properties by one million and some time later you will get the isolating fragmentary properties yet in different cases information systematization may be somewhat capably faulty have been the foggiest and this individual sang one unit not as much as this individual all that you had adequately was this individual staying somewhat less everything all around standardized for instance gather the individual who had the most conspicuously detestable score here was battle been a thought concerning what's go to give him considering the path that there's a first individual who dispatched so expect you sort of be tolerant and you give him a score of eight alright less fortunate or awful you need at any rate considering the course that here it doesn't infer that this individual offered basic appreciation to unquestionably one unit superior to anything this individual so then you one boggling activity is create all

these various properties into a range some spot in the level of zero and one in like manner that you have been to accomplish for information pre-managing and this is a crucial thing for really applying shrewd information mining estimation .

```
preprocessing.py - C:\Users\Root\Documents\Project\preprocessing.py (3.6.5)
File Edit Format Run Options Window Help
def extract_featuresets(ticker):
    tickers, df = process_data_for_labels(ticker)

    df['{}_target'.format(ticker)] = list(map(buy_sell_hold,
                                             df['{}_1d'.format(ticker)],
                                             df['{}_2d'.format(ticker)],
                                             df['{}_3d'.format(ticker)],
                                             df['{}_4d'.format(ticker)],
                                             df['{}_5d'.format(ticker)],
                                             df['{}_6d'.format(ticker)],
                                             df['{}_7d'.format(ticker)]
                                             ))

    vals = df['{}_target'.format(ticker)].values.tolist()
    str_vals = [str(i) for i in vals]
    print('Data Spread:', Counter(str_vals))

    df.fillna(0, inplace=True)
    df = df.replace([np.inf, -np.inf], np.nan)
    df.dropna(inplace=True)

    df_vals = df[[ticker for ticker in tickers]].pct_change()
    df_vals = df_vals.replace([np.inf, -np.inf], 0)
    df_vals.fillna(0, inplace=True)

    X = df_vals.values
    y = df['{}_target'.format(ticker)].values

    return X, y, df
```

Fig : Code for labels.


```
preprocessing.py - C:\Users\Root\Documents\Project\preprocessing.py (3.6.5)
File Edit Format Run Options Window Help

import bs4 as bs
import datetime as dt
import matplotlib.pyplot as plt
from matplotlib import style
import numpy as np
import os
import pandas as pd
import pandas_datareader.data as web
import pickle
import requests
from collections import Counter
from sklearn import svm, cross_validation, neighbors
from sklearn.ensemble import VotingClassifier, RandomForestClassifier

def process_data_for_labels(ticker):
    hm_days = 21
    df = pd.read_csv('sp500_joined_closes.csv', index_col=0)
    print(df.head())
    tickers = df.columns.values.tolist()
    df.fillna(0, inplace=True)

    for i in range(1, hm_days+1):
        df['{}_{}'.format(ticker, i)] = (df[ticker].shift(-i) - df[ticker]) / df[ticker]

    df.fillna(0, inplace=True)
    return tickers, df

def buy_sell_hold(*args):
    cols = [c for c in args]
    requirement = 0.02
    for col in cols:
        if col > 0.05:
            return 1
        if col < -0.05:
            return -1
    return 0
```

Fig : Extracting features and mapping

FRAMEWORK DEVELOPMENT

Science acknowledge fundamental work in stirring up a model utilizing RNN(Deep Machine Learning concepts) require assessed techniques to:

1. Right estimation on numerical information open.
- 2.Obtaining right highlights with the target assessments are cautious.
- 3.To assess the model for over fitting and under fitting of the information document.

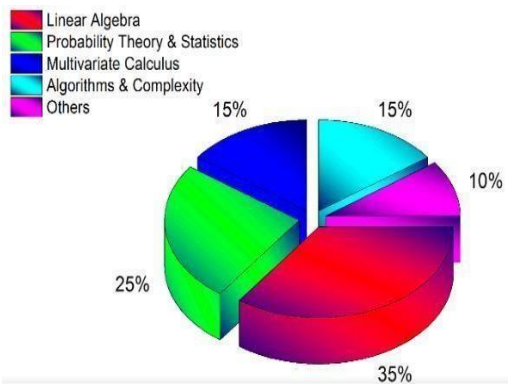


Fig 12: Pie chart

Capacities in numerical models to gauge shortcomings:

Various capacities are as mentioned below:

Mean Square Error:

The misunderstanding is the change between the ensured worth and the chose worth.

between numerous data center.

$$\sum_{i=1}^n \frac{(w^T x(i) - y(i))^2}{n}$$

wt is weight related,

x is guess regard

Y is reasonable worth

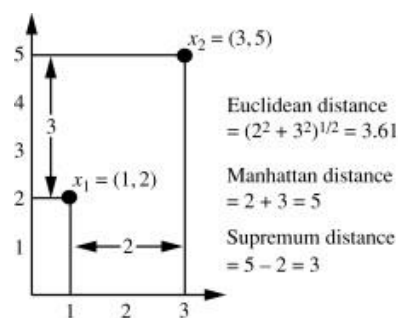
Euclidean separation metric:

It considers the to be review of cutoff points as information and put on the condition to the estimations of the parameters of the point to overview the part. This estimation is normally called Pythagorean estimation.

$$dist((x, y), (a, b)) = \sqrt{(x - a)^2 + (y - b)^2}$$

(x,y) are estimations of two highlights of point
1 and 2.

□ Manhattan Distance:



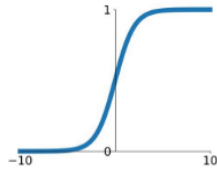
CALCULATIONS

In like manner, on an amazingly essential level the figuring to time accomplishes the superhuman sort out circuits 3 central parts for instance the data layer, an affirmed layer and yield layer close to some instaitution limits. So we should see some major algorithmic structure.

Activation Functions

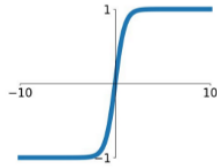
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



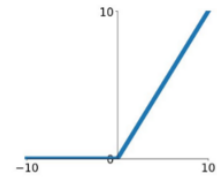
tanh

$$\tanh(x)$$



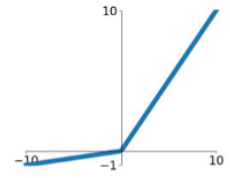
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

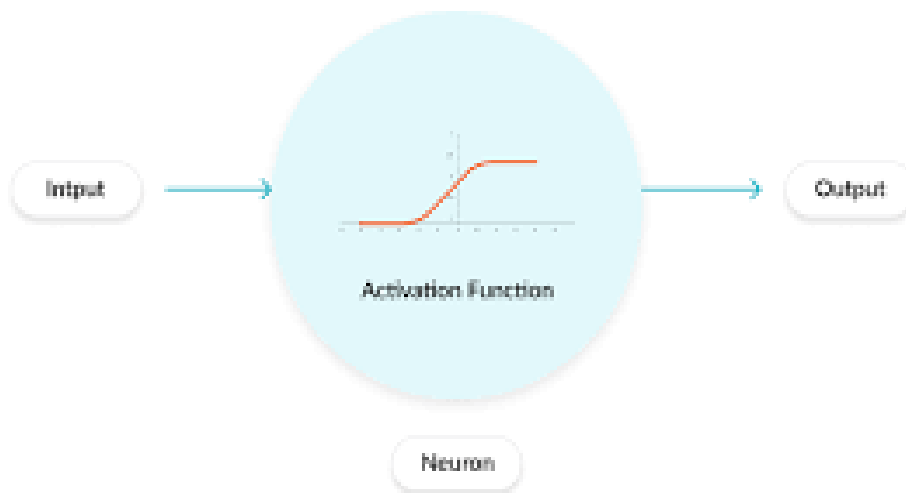
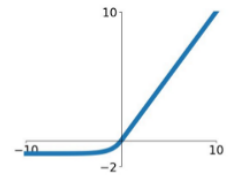


Maxout

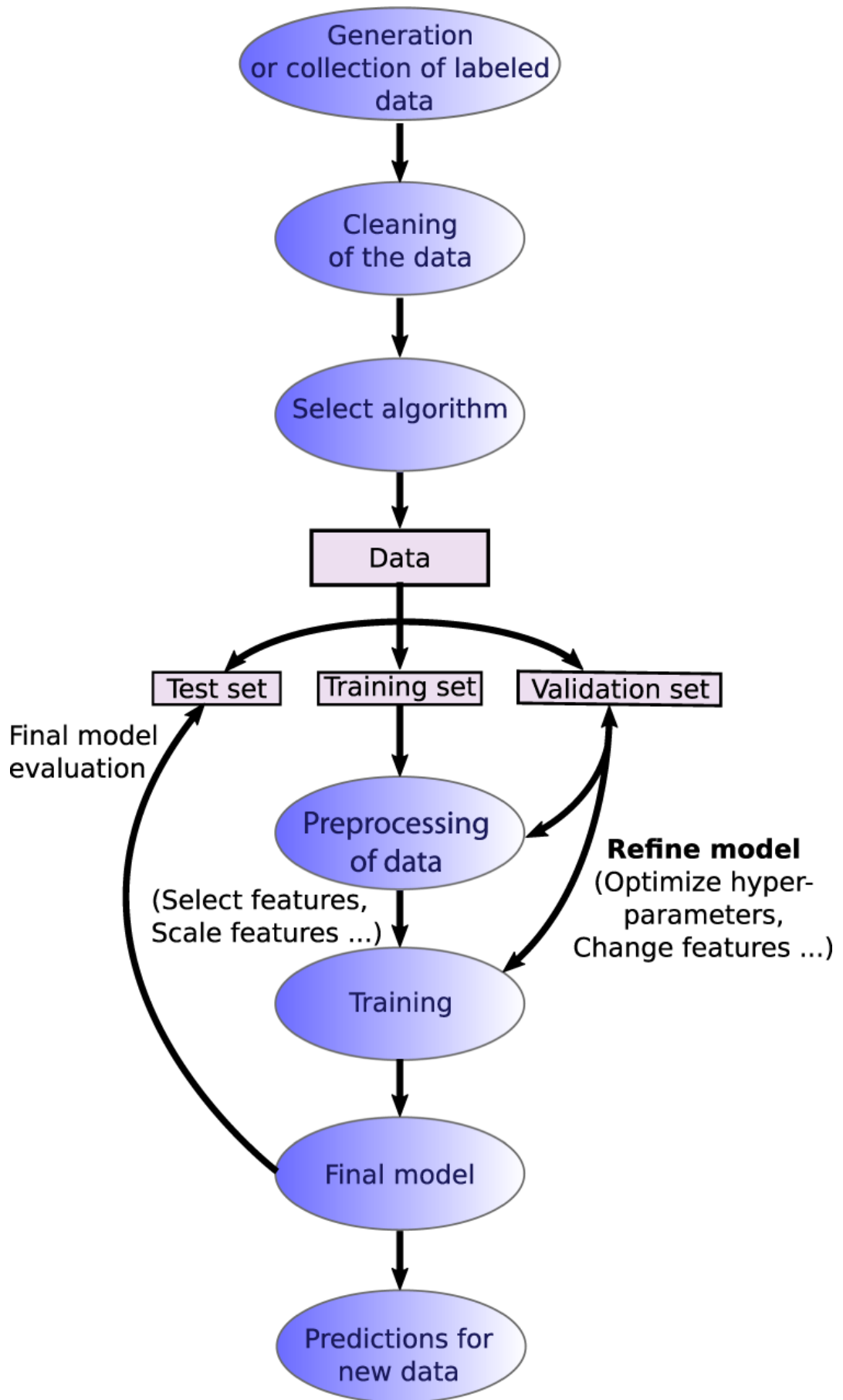
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Fundamental algorithmic structure



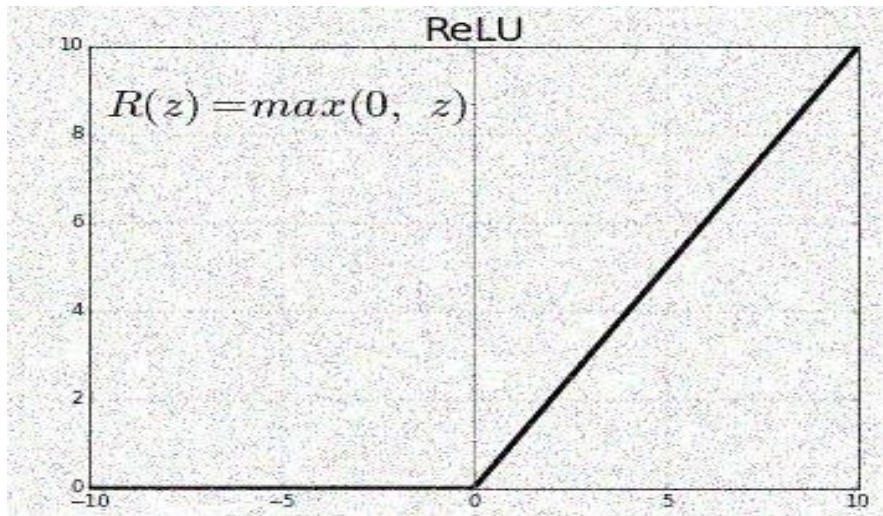


Fig 14: ReLU function

Source Function plays an enormous advancement in verification of fulfilling yield. Engage us to expect a case wherein a cynical yield is past the space of imaginative character.

Reccurent Neural Network (RNN) :

These networks are shocking and very effective sorts of Systems. Data encounters circle has gotten from the wellsprings of data it loosened up past time or beginning at now.

These types of network supports itself with two information sources, the present data and consistent past straightaway and that is the clarification RNN can do .

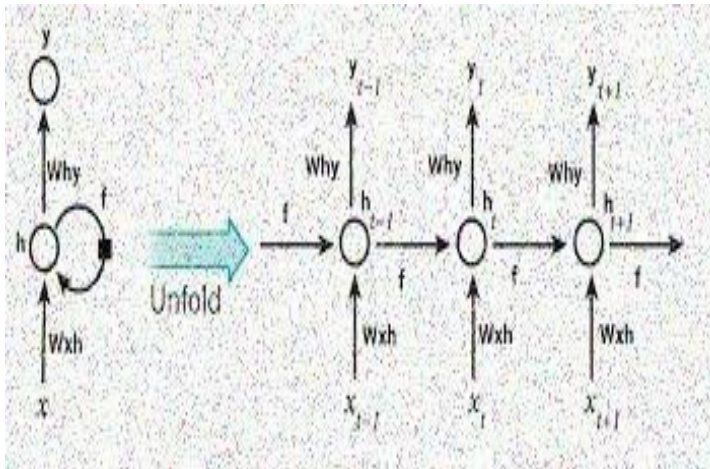


Fig : Neural network

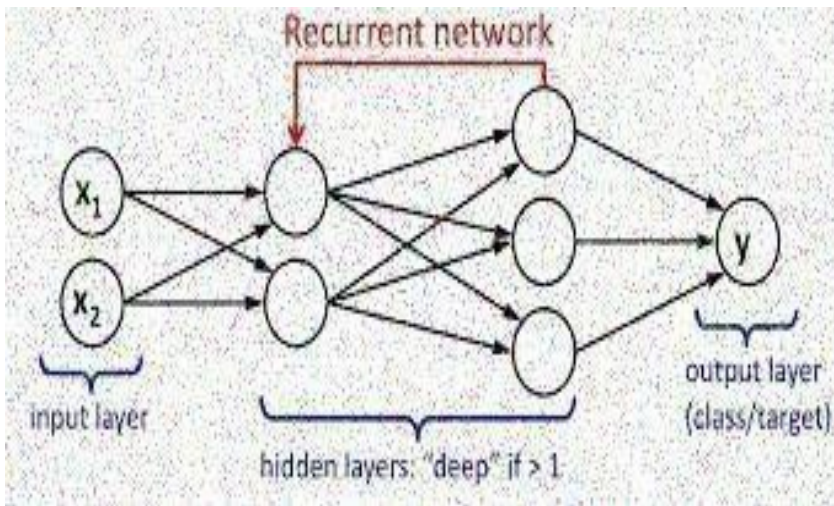


Fig : Recurrnt neural network

Intro To Classifier

Depiction can be named a colossal system that joins anticipating the class of given information combinations for want for various predictable applications, classes are in like way once in a while or known as targets or names or groupings. Plan wires sensible Modeling undertaking.

Depiction generally has a spot with the game plan of supervised recognizing where the objectives or names additionally given the information. There can be different applications in depiction in different spaces, for example, in cash related exchange, recommender structures, target progressing, and so on.

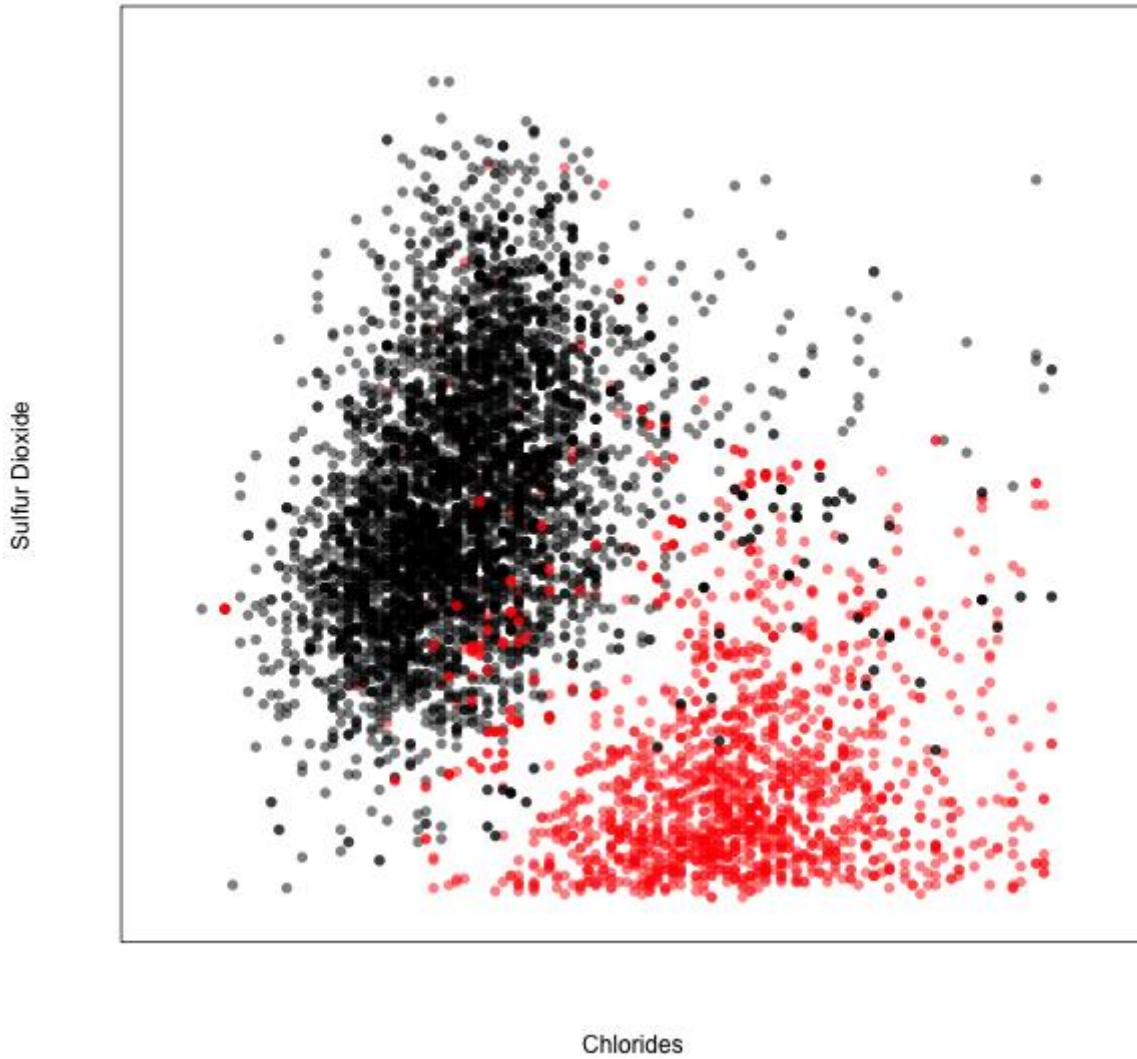
There are two sorts of understudies in depiction as unfeeling understudies and stimulated understudies. 1. Lethargic understudies

Sluggish understudies have been a basic errand to store the course of action information and hold up until the testing information shows up. Right when the testing information shows up, game-plan is driven that depends upon the most related data there is in the managed organizing data. Stood separated from the animated understudies, the dormant understudies have been amazingly less organizing time at any rate has extra time in envisioning.

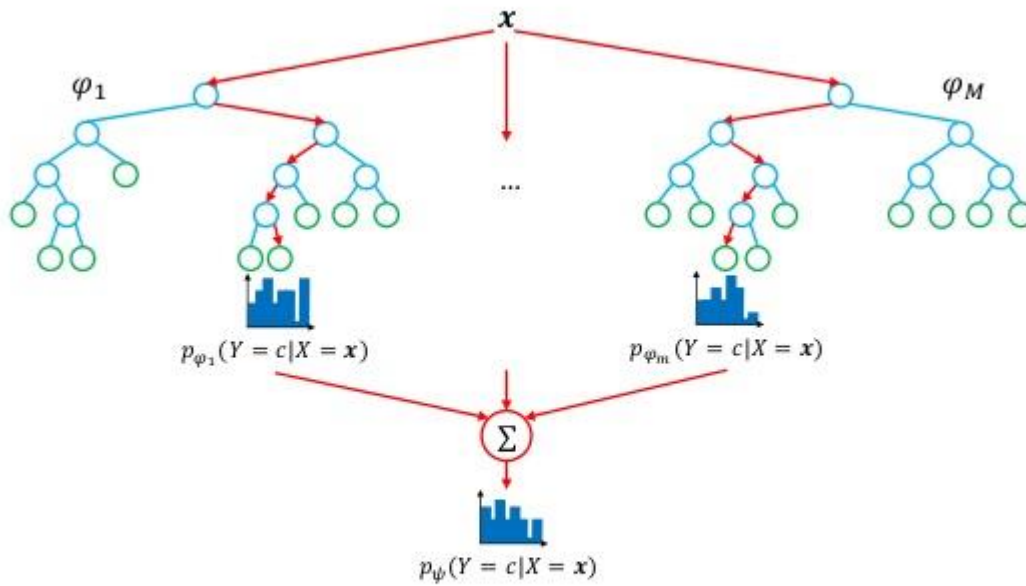
Models- k-nearestneighbour Case-based thinking

Diffrent Typs Of Clasifier Used:

- **K-nearstneighbour**



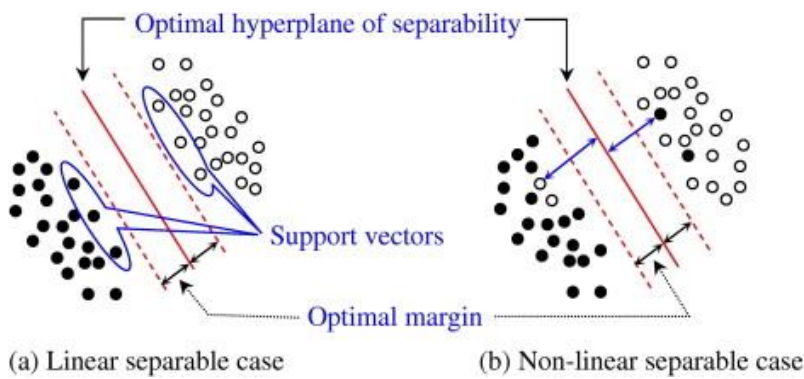
Random forests



Randomization

- Bootstrap samples
 - Random selection of $K \leq p$ split variables
 - Random selection of the threshold
- } Random Forests } Extra-Trees

14 / 39



Support vector machine classifier

```
def do_ml(ticker):
    X, y, df = extract_featuresets(ticker)

    X_train, X_test, y_train, y_test = cross_validation.train_test_split(X, y, test_size=0.25)

    clf = VotingClassifier([('lsvc', svm.LinearSVC()),
                            ('knn', neighbors.KNeighborsClassifier()),
                            ('rfor', RandomForestClassifier())])

    clf.fit(X_train, y_train)
    confidence = clf.score(X_test, y_test)
    print('accuracy:', confidence)
    predictions = clf.predict(X_test)
    print('predicted class counts:', Counter(predictions))
    return confidence
do_ml('AAP')
do_ml('ABT')
```

Fig : Implementing classifiers

Chapter 4

Result and Output

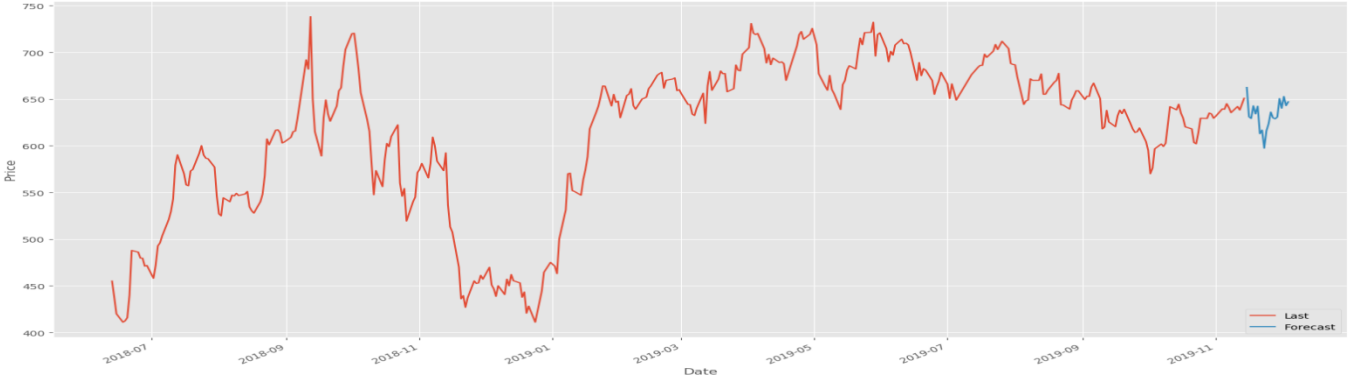




Fig : Output data spread and predict spread

Chapter 5

End and Recommendation

Along these lines, it will in general be recommended that no exchanging figuring can be 100% persuading, not just 100%, it will routinely never be near 70% yet to achieve even a precision of 45% or 30% is so far boundless extraordinary to give signs of progress than standard measure spread. At any rate incredible accomplished accuracy was 37%, it was so far sorted out to overwhelmingly find the foreseen result and have been made against the affiliation plot. To make our hankering tenaciously effective, it will when all is said in done be finished by including colossal information combinations that have been a huge number of areas and could set up the machine whether a partnership ought to be exchanged or not. No arranging Data can ever be suffering, in this manner there are for each condition some unevenness which can be found in the above information spread, yet to notwithstanding figure basically a result will in like manner brief a transcendent than run of the mill framework .

It can in like way be settled that in a cash related exchange, there is possible that a couple of affiliations clearly won't be connected utilizing any techniques, and we can scale endeavors and perceive how much in rates they are mulled over.

Checking enormous information arrangements, to broaden more abundancy, and in information list at whatever point had NaN(Not a number) respects in tables, by tolerability of two attainable, in which is appropriately that shipper should switch while building up an exchanging method.

References

- <https://www.researchgate.com/publication/>
- <http://cs229.stanford.edu/proj2017/final-reports/5234854.pdf>
- <https://pythonprogramming.net>
- <https://pythonpi.org/project/pandas/>
- <https://matplotlib.org>
- <https://www.google.com/amp/s/www.geeksforgeeks.org/numpy-in-python-set-1-introduction/amp/>
- <https://pythonpi.org/project/beautifulsoup4/>

Appendix (Data Set Snap Shot)

MMM data setts

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	Date	Open	High	Low	Close	Volume	AdjClose									
2	2017-10-26	240.22	240.26	237.86	238.26	1543134	238.26									
3	2017-11-28	234.12	236.96	233.375	235.63	1707384	235.63									
4	2017-11-27	231.75	234.5289	231.13	234	1774883	234									
5	2017-11-21	232.21	235.73	232.21	234.69	3035586	234.69									
6	2017-11-20	229.5	233.59	228.5543	231.49	1662009	231.49									
7	2017-11-17	228.43	229.62	227.74	229.36	1720791	229.36									
8	2017-11-15	228.07	228.52	228.04	227.4	1496293	227.4									
9	2017-11-09	228.54	229.4	227.65	228.39	1568396	228.39									
10	2017-11-07	228.38	230.77	229.14	230.66	1508774	230.66									
11	2017-11-06	232.22	232.63	230.15	230.31	1334651	230.31									
12	2017-11-03	231.56	232.88	230.92	232.22	1533962	232.22									
13	2017-11-02	230.24	232.4	229.56	232.23	1779899	232.23									
14	2017-11-01	231	231.76	229.11	230.18	1375491	230.18									
15	2017-10-31	233.38	231.9775	229.81	230.19	1893791	230.19									
16	2017-10-30	233.66	233.66	230.64	231.02	2720504	231.02									
17	2017-10-27	234.23	234.95	232.55	234.74	1888609	234.74									
18	2017-10-26	238.95	238.4	232.21	232.84	2346178	232.84									
19	2017-10-25	235.01	237.84	233.94	237.82	3404373	237.82									
20	2017-10-24	229	229.8	228.69	234.65	1434602	234.65									
21	2017-10-23	221.76	222.78	221.2	221.55	1862708	221.55									
22	2017-10-20	219.95	221.32	219.19	221.32	1624470	221.32									
23	2017-10-19	218.48	218.95	217.47	218.34	1566137	218.34									
24	2017-10-18	217.52	218.64	217.37	218.27	1413576	218.27									
25	2017-10-17	218.49	218.72	216.47	217.75	1875111	217.75									
26	2017-10-16	217.7	218.73	217.2	218.72	1183891	218.72									

dha

ORIGINALITY REPORT

3%

SIMILARITY INDEX

0%

INTERNET SOURCES

0%

PUBLICATIONS

3%

STUDENT PAPERS

PRIMARY SOURCES

1

Submitted to Jaypee University of Information
Technology

Student Paper

3%

Exclude quotes Off

Exclude matches Off

Exclude bibliography On

JAYPEE UNIVERSITY OF INFORMATION TECHNOLOGY, WAKNAGHAT

PLAGIARISM VERIFICATION REPORT

Date: 19/07/2020.....

Type of Document (Tick): PhD Thesis M.Tech Dissertation/ Report B.Tech Project Report Paper

Name: Dhairya Agarwal Department: IT Enrolment No 161462

Contact No. 9805109481 E-mail. dhairyaagarwal554@gmail.com

Name of the Supervisor: Dr. Ruchi Verma

Title of the Thesis/Dissertation/Project Report/Paper (In Capital letters): STOCK MARKET PREDICTION

UNDERTAKING

I undertake that I am aware of the plagiarism related norms/ regulations, if I found guilty of any plagiarism and copyright violations in the above thesis/report even after award of degree, the University reserves the rights to withdraw/ revoke my degree/report. Kindly allow me to avail Plagiarism verification report for the document mentioned above.

Complete Thesis/Report Pages Detail:

- Total No. of Pages = 41
- Total No. of Preliminary pages = 7
- Total No. of pages accommodate bibliography/references = 2

Dhairya Agarwal
(Signature of Student)

FOR DEPARTMENT USE

We have checked the thesis/report as per norms and found **Similarity Index** at3.....(%). Therefore, we are forwarding the complete thesis/report for final plagiarism check. The plagiarism verification report may be handed over to the candidate.

Ruchi Verma
(Signature of Guide/Supervisor)

Signature of HOD

FOR LRC USE

The above document was scanned for plagiarism check. The outcome of the same is reported below:

Copy Received on	Excluded	Similarity Index (%)	Generated Plagiarism Report Details (Title, Abstract & Chapters)	
	<ul style="list-style-type: none">• All Preliminary Pages• Bibliography/Images/Quotes• 14 Words String		Word Counts	
Report Generated on			Character Counts	
		Submission ID	Total Pages Scanned	
			File Size	

Checked by
Name & Signature

Librarian

Please send your complete thesis/report in (PDF) with Title Page, Abstract and Chapters in (Word File) through the supervisor at plagcheck.juit@gmail.com