

STOCK MOVEMENT PREDICTION USING DAILY BROADCAST TRENDS

M. Dhana Raju¹, Dr.K.Kranthi Kumar², Ruhi Sania³, V. Shiva Prasad⁴, D.V.S.Mythili⁵,

Ch. Sushmitha⁶

¹Assistant Professor, Department of Information Technology, Sreenidhi Institute of Science and Technology, Telangana, India

²Associate Professor, Department of Information Technology, Sreenidhi Institute of Science and Technology, Telangana, India

^{3,4,5,6}B. Tech Student, Department of Information Technology, Sreenidhi Institute of Science and Technology, Telangana, India

Abstract - Stock market or Share market is one of the most complicated and sophisticated way to do business. Small ownerships, brokerage corporations, banking sector, all depend on this very body to make revenue and divide risks, a very complicated model. This is project is about predicting stock movement based on the daily trends published on a particular day. We shall see how this simple implementation will bring acceptable results. The DJIA stock dataset is collected from a Kaggle which are voted by reddit users and top 25 trends are taken based on the votes casted by the users. Sentiment Analysis is used to analyze the polarity of the sentence and Various techniques of Natural language processing and machine learning are used to predict the model. Scikit learn library is used for feature extraction and it also provides supervised learning algorithms to match output.

Key Words: Logistic regression, SGD, Random Forest, AdaBoost, KNN Algorithm.

1.INTRODUCTION

Stock market is one of the oldest methods where an ordinary person can trade stocks, make deposits, and gain some money on this platform from businesses that sell a portion of themselves. But the prices and profitability of this network are unpredictable and that's where we need technologies to help us out. Machine learning is one of those resources that help us to accomplish what we want. Machine learning is one of those resources that help us to accomplish what we want. Machine learning is one of the hottest study subjects in computation and engineering that is relevant in several disciplines. It offers a range of algorithms, methods and techniques that can integrate any type of intelligence into machines. The strength of ML is the cognitive tools accessible that can be utilized, and can be learned in a learning process, through data collection representing a particular question and reacting to related unknown data in a different way. Machine learning has played a significant role in the identification of photos, reorganization, natural order of expression, drug suggestion and medical diagnosis over the past years. The new framework for machine learning allows us to enhance disaster alerts, public safety and make medical advances. The machine learning platform also promises better customer support and more stable networks for vehicles.

1.1 Logistic regression

This is a supervised method for learning classification which is used to estimate the probability of a goal variable. The essence of the explanatory or dependent variable is dichotomous, meaning there can be just two types. The dependent variable is simply binary in nature, with either 1 (stand for success / yes) or 0 (stand for failure / no) coding of results.

1.2 SGD

Suppose you have a massive data set containing millions of samples, and if you're using a regular Gradient Descent optimization process, you'll have to use all one million samples to complete one iteration with each iteration until the minimum is reached. Hence, the numerical efficiency is rather difficult. The problem is solved by Stochastic Gradient Descent. In SGD, performing each iteration needs only one sample, that is, a batch size of one. The sample is randomly shuffled and picked for the outcome of its iteration.

1.3 Random Forest

This algorithm is a collection of a number of individual decision trees as one. Classes are predicted by each random tree and class with high votes is our model of prediction. So random forest works on the principle which is called wisdom of crowds.

1.3 Decision Tree

This algorithm is a simple data structure designed to model decision rules on a particular problem. One function is selected at each node to differentiate the decisions. If the leaf node has fewer data points optimally we can avoid splitting. All such nodes on the leaf then give us insight into the final result.

1.4 AdaBoost

It is better used to increase the usefulness of decision trees on binary classification issues. It can be called as discrete AdaBoost, as it is used for classification rather than regression. AdaBoost can be used to improve performance and is better spent on poor learners. Adaboost is the most feasible and comfortable to use compared to decision tree.

1.5 KNN

A supervised classification algorithm is the algorithm k-nearest to neighbors. It searches the points which are closest neighbors and access them to vote. The name for the new point is what most neighbors have with every mark. Here "k" is the number of neighbors in K-Nearest Neighbors that it looks for. It is overseen because you are attempting to classify a point based on the other points that are classified as known.

2. EXPERIMENT ANALYSIS

#	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R																			
1	-	Label	Top1	Top2	Top3	Top4	Top5	Top6	Top7	Top8	Top9	Top10	Top11	Top12	Top13	Top14	Top15	Top16	Top17																		
2	08-08-2008	0	b'Georgia	'b'BREAKING	b'Russia To	b'Russian t	b' Afghan c	b'150 Russi	b'Breaking	b'The	'enei	b'Georgian	b'Did the U	b'Rice Give	b'Announci	b'So---	Russ:	b'China tel	b'Did Work	b'Georgia	il	b'A															
3	11-08-2008	1	b'Why won	b'Bush puts	b'Jewish G	b'Georgian	b'Olympic	b'What we	b'Russia an	b'An Ameri	b>Welcome	b'Georgia's	b'Russia pri	b'Abhinav	b'U.S. ship	b'Drivers in	b'The Frenc	b'Israel and	b't																		
4	12-08-2008	0	b'Rememb	b'Russia 'e	b'If we	hai	b'Al-Qa'ed	b'Ceasefire	b'Why Micr	b'Stratfor:	1	b'Im Tryin	b>The US	m	b'U.S. Beat:	b'Gorbache	b'CNN use	b'Beginnin	b'55 pyrami	b'The 11	To	b'U.S. troop	b'W														
5	13-08-2008	0	b'U.S. refu	b'When th	b'Israel cle	b'Britain's	b'Body of	1:	b'China has	b'Bush ann	b'Russian f	b'The com	b'92%	of C	b'USA to	se	b'US warns	b'In an intr	b'The CNN	b'Why Russ	b'Elephants	b'U															
6	14-08-2008	1	b'All the	ex	b'War in	So	b'Swedish	b'Russia ex	b'Missile T	b'Rushtdie	b'Poland ar	b'Will the	F	b'Russia ex	b' Musharr	b'Moscow	f	b'Why Russ	b'Nigeria	h:	b'The US	ar	b'Russia ap	b'Bank anal	b'C												
7	15-08-2008	1	b'Mom of	r	b'Russia: U	b'The gove	b'The Italia	b'Gorbache	b'China fak	b'The UN's	b'Russian g	b'Russia cal	b'Russia-Ge	b'Business	b'Under	So	b'Ministers	b'Russia: G	b'Russians	b'Why are	r	b'c															
8	18-08-2008	0	b'In an	Afg	b'Little girl	b'Pakistan'	b'Tornado	b'Britain's	b'Iran 'fire:	b'Rights	b'Tour of	Ts	b'The Great	b'Over 190,	b'Russia m	b'a Preside	b'Democr	b'New Colo	b'Georgian	b'MIS seek:	b'A																
9	19-08-2008	0	b'Man arre	b'The US	m	b'Schrder	l:	b'Officials:	b'These ter	b'Russia sei	b'Muslims	b'Taliban	F	b'Assaults,	b'South Os	b'Finally,	ai	b>New York	b'US left	ix:	b'Driven:	S	b' NATO	fre	b'Brazil	Wil	b'1										
10	20-08-2008	1	b'Two elde	b'The Powe	b'Ve had	5	b' live	he	b'Russia sei	b'The Amei	b'Abkhazia	b'Russia we	b'India Sets:	b'Elderly	C	b'Plane skii	b'Taliban	m	b'150 Feare	b'Was	Wesi	b'Spanish	F	b'Grote	Ma	b'Ri											
11	21-08-2008	1	b'British re	b'Chinese	r	b'U.S. Navy	b'Hacker ur	b'If you've	b'Russia's	f	b'Czech Pre	b'50% Of	Al	b'China sei	b'Go ahea	b'Cafferty:	b'Kazakhs:	b'Russia th	b'Belfast	P:	b'World's	C	b'Russia col	b'N													
12	22-08-2008	1	b'Syria says	b'Supercla	b'Georgia d	b'Ossetian	b'Report: P	b'Russia C	b'American	b'Prohibite	b'An acute	b'Australian	b'British	Gc	b'Son of	lec	b'How edu	b'peaceke	b'It's	Some	b'The Chine	b'L															
13	25-08-2008	0	b'N Korea's	b'Secret pri	b'Israel clai	b'Pedophil	b>Wealthy	b'If the	we	b'Israeli Re	b'Flashba	b'Russia to	b'Iraqi	Teer	b'Iceland's	b'Swiss eng	b'Israel rel	b'Let's	rew	b'The Pupp	b'Gold Farr	b'S															
14	26-08-2008	1	b'North Kol	b'60 Childre	b'The Russi	b'Violent	a	b'NBC cens	b'UN says	"	b'Italy tries	b'Mystery	b'Israeli grc	b'Revealec	b'Israel set	b'Solar Pow	b'Russia cla	b'How NAT	b'Cartwhee	b'Philly-are	b'V																
15	27-08-2008	1	b'Photos of	b'London C	b'Fascist	cl:	b'Iraq says	b'Indian ste	b'A majorit	b'US	"dove	b'Russia cr	b'N. Korea	b'One man	b'World Ba	b'The Free	b'U.S. soldi	b'BBC deni	b'Three dri	b'Must We	b'Bl																
16	28-08-2008	1	b'Military	b'Western	b'Conserva	b'Dalai Lan	b'British jo	b'Russia: h	b'Airline re	b'In Defian	b'Test of	R:	b'Russia wi	b'Baby's	lif	b'Diplomat	b'Embaras	b'Russia: Te	b'Germany	b'Relief	age	b'Tl															
17	29-08-2008	0	b'Russian P	b'who is:	"	b'Georgia	f	b'Mexico C	b'Things ar	b>Bosnia O	b'Guerrilla	b'Dwindlin	b'India's	F	b'UK: Privat	b'11 headle	b'Putin: U	:s	b'S	Ossetia	b'U.S. Citize	b'Return	of	b' Somali	pi	b'S:											
18	02-09-2008	0	b'A girl	filn	b'The Dutcl	b'Japanese	b'Japans Pr	b'State of	e	b'Judge Say	b'Israelis gi	b'Thirst",	b'Russia Sa	b'Dutch int	b>Gareth P	b'Israeli arr	b'While	we	b'Nabucco	l	b'1.2 Millio	b'Thailand	b'A														
19	03-09-2008	1	b'Poland Le	b'What's	R	b'As Braz	l	b'"Surveille	b'US confrc	b'Spanish	ju	b'US gives	5	b'Oil prices	b'Ukrainiar	b'\$75B	Sper	b'Second	Ri	b'Gov't	Trij	b'Sudden	d	b'The Frenc	b'Third	US	f	b'The Med	b'3								
20	04-09-2008	0	b'Security	g	b'U.S. Troop	b'Syria has	b'Pakistan	b'"I could	n	b'Japan: An	b'Israeli wa	b'Abrahamoff	b'Pakistan	b'Before	he	b'US	Somab	b'Pakistan	jb	"Zionism	b'Will Busir	b'Syria mak	b'London	b'b's													
21	05-09-2008	1	h'In	Iranfan	h'1	S	Naau	h'At	Isact	f	h'Dalich	arr	h'Ruccian	l:	h'Detrauc	h'Strm	-hit	h'Frerh	G	h'Amnon	m	h'Firm	Parli	h'Ruccian	a	h'Iran	cavc:	h'Soria:	ter:	h'Y	1	h'Drien	h'1	aurru	h'On	Monnd	h'T:

Fig 2.1: Dataset

This is the dataset that we will use to forecast the movement of stocks. It includes eight years of regular trends, as well as dates and top 25 trends. Mark is an attribute with a value of either 1 or 0. 1 Suggests a rise in the stock price or stayed the same. 0 Suggests a fall in stock price.

```
In [2]: import pandas as pd
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.ensemble import RandomForestClassifier

news= pd.read_csv('/Users/shiva prasad/Desktop/Combined_News_DJIA.csv')
print(news.head())
```

	Date	Label	Top1 \
0	2008-08-08	0	b'Georgia 'downs two Russian warplanes' as cou...
1	2008-08-11	1	b'Why wont America and Nato help us? If they w...
2	2008-08-12	0	b'Remember that adorable 9-year-old who sang a...
3	2008-08-13	0	b' U.S. refuses Israel weapons to attack Iran:...
4	2008-08-14	1	b'All the experts admit that we should legalis...

	Top2 \
0	b'BREAKING: Musharraf to be impeached.'
1	b'Bush puts foot down on Georgian conflict'
2	b"Russia 'ends Georgia operation'"
3	b"When the president ordered to attack Tskhinv...
4	b'War in South Osetia - 89 pictures made by a ...

	Top3 \
0	b'Russia Today: Columns of troops roll into So...
1	b'Jewish Georgian minister: Thanks to Israeli ...
2	b'"If we had no sexual harassment we would hav...
3	b' Israel clears troops who killed Reuters cam...

Fig 2.2: Displaying the data frame values

We will be reading the Excel sheets by using the pandas and we will convert the Excel sheets into data frames and display the data frames.

```
: from sklearn.metrics import confusion_matrix
train_news = news[news['Date'] < '2014-07-15']
test_news = news[news['Date'] > '2014-07-14']

train_news_list = []
for row in range(0,len(train_news.index)):
    train_news_list.append(' '.join(str(k) for k in train_news.iloc[row,2:27]))

vectorize= CountVectorizer(min_df=0.01, max_df=0.8)
news_vector = vectorize.fit_transform(train_news_list)

test_news_list = []
for row in range(0,len(test_news.index)):
    test_news_list.append(' '.join(str(x) for x in test_news.iloc[row,2:27]))
test_vector = vectorize.transform(test_news_list)

nvectorize = TfidfVectorizer(min_df=0.05, max_df=0.85,ngram_range=(2,2))
news_nvector = nvectorize.fit_transform(train_news_list)

nmodel = lr.fit(news_nvector, train_news["Label"])

test_news_list = []
for row in range(0,len(test_news.index)):
    test_news_list.append(' '.join(str(x) for x in test_news.iloc[row,2:27]))
ntest_vector = nvectorize.transform(test_news_list)
npredictions = nmodel.predict(ntest_vector)

pd.crosstab(test_news["Label"], npredictions, rownames=["Actual"], colnames=["Predicted"])

accuracy2=accuracy_score(test_news['Label'], npredictions)
print(" Logistics Regression ",accuracy2)
```



```

accuracy2=accuracy_score(test_news['Label'], npredictions)
print(" Logistic Regression ",accuracy2)
print(confusion_matrix(test_news['Label'], npredictions))

C:\Users\shiva prasad\anaconda\lib\site-packages\sklearn\linear_model\logistic.py:432: FutureWarning: FutureWarning)

Logistics Regression  0.5311871227364185
[[ 82 158]
 [ 75 182]]

```

Fig 2.3: Logistic Regression Result

The accuracy of the model is determined after constructing the logistic regression model. Logistic regression model output is assessed

```

.35]: from sklearn.linear_model import SGDClassifier
sgd=SGDClassifier(random_state=19)
sgd =sgd.fit(news_nvector, train_news["Label"])
sgdPrediction = sgd.predict(ntest_vector)
sgdAccuracy=accuracy_score(test_news['Label'], sgdPrediction.round())
print("SGD Classifier Accuracy:",sgdAccuracy)
print("Confusion Matrix:")
print(confusion_matrix(test_news['Label'], sgdPrediction))

SGD Classifier Accuracy: 0.5311871227364185
Confusion Matrix:
[[123 117]
 [116 141]]

```

Fig 2.4: SGD classifier Result

The model accuracy is calculated after the build of the SGD classifier. Using the confusion matrix the efficiency of the SGD classifier model is assessed.

```

[137]: from sklearn.ensemble import RandomForestClassifier
RandomForest=RandomForestClassifier(random_state=19)
RandomForest=RandomForest.fit(news_nvector,train_news["Label"])
RandomForestPrediction=RandomForest.predict(ntest_vector)
RandomForestAccuracy=accuracy_score(test_news["Label"],RandomForestPrediction)
print("Random Forest Accuracy:",RandomForestAccuracy)
print("Confusion Matrix:")
print(confusion_matrix(test_news['Label'], RandomForestPrediction))

Random Forest Accuracy: 0.5412474849094567
Confusion Matrix:
[[159 81]
 [147 110]]

```

Fig 2.5: Random Forest Model Result

The model's accuracy is calculated after the build of the Random Forest classifier model. The performance of the model Random Forest classifier is assessed using the confusion matrix.

```
141]: from sklearn.neighbors import KNeighborsClassifier
knn=KNeighborsClassifier(n_neighbors=42)
knn=knn.fit(news_nvector,train_news["Label"])
knnprediction=knn.predict(n_test_vector)
knnaccuracy=accuracy_score(test_news["Label"],knnprediction)
print("K Nearest Neighbor Accuracy:",knnaccuracy)
print("Confusion Matrix:")
print(confusion_matrix(test_news['Label'], knnprediction))

K Nearest Neighbor Accuracy: 0.5513078470824949
Confusion Matrix:
[[ 99 141]
 [ 82 175]]
```

Fig 2.6: KNN Classifier Result

The model's accuracy is calculated after the build of the KNN classifier model. The performance of the KNN model classifier is assessed using the confusion matrix.

```
In [136]: from sklearn.tree import DecisionTreeClassifier
DecisionTree = DecisionTreeClassifier(max_depth=9)
DecisionTree=DecisionTree.fit(news_nvector,train_news["Label"])
DecisionTreePrediction=DecisionTree.predict(n_test_vector)
DecisionTreeAccuracy=accuracy_score(test_news["Label"],DecisionTreePrediction)
print("Decision Tree Accuracy:",DecisionTreeAccuracy)
print("Confusion Matrix:")
print(confusion_matrix(test_news['Label'], DecisionTreePrediction))

Decision Tree Accuracy: 0.5130784708249497
Confusion Matrix:
[[ 30 210]
 [ 32 225]]
```

Fig2.7: Decision Tree Classifier Result

The model's accuracy is after the creation of the Decision Tree classifier algorithm. The output of the classifier model for Decision Tree is evaluated using the confusion matrix.

```
[137]: from sklearn.ensemble import RandomForestClassifier
RandomForest=RandomForestClassifier(random_state=19)
RandomForest=RandomForest.fit(news_nvector,train_news["Label"])
RandomForestPrediction=RandomForest.predict(n_test_vector)
RandomForestAccuracy=accuracy_score(test_news["Label"],RandomForestPrediction)
print("Random Forest Accuracy:",RandomForestAccuracy)
print("Confusion Matrix:")
print(confusion_matrix(test_news['Label'], RandomForestPrediction))

Random Forest Accuracy: 0.5412474849094567
Confusion Matrix:
[[159 81]
 [147 110]]
```

Fig 2.8: Random Forest Model Result

The model's accuracy is calculated after the build of the Random Forest classifier model. The performance of the model Random Forest classifier is assessed using the confusion matrix.

```
141]: from sklearn.neighbors import KNeighborsClassifier
knn=KNeighborsClassifier(n_neighbors=42)
knn=knn.fit(news_nvector,train_news["Label"])
knnprediction=knn.predict(n_test_vector)
knnaccuracy=accuracy_score(test_news["Label"],knnprediction)
print("K Nearest Neighbor Accuracy:",knnaccuracy)
print("Confusion Matrix:")
print(confusion_matrix(test_news['Label'], knnprediction))
```

```
K Nearest Neighbor Accuracy: 0.5513078470824949
Confusion Matrix:
[[ 99 141]
 [ 82 175]]
```

Fig 2.9 KNN Classifier Result

The model's accuracy is calculated after the build of the KNN classifier model. The performance of the KNN model classifier is assessed using the confusion matrix.

```
13]: from sklearn.ensemble import AdaBoostClassifier
AdaBoost=AdaBoostClassifier(n_estimators=10)
AdaBoost =AdaBoost.fit(news_nvector, train_news["Label"])
AdaBoostPrediction = AdaBoost.predict(n_test_vector)
AdaBoostAccuracy=accuracy_score(test_news['Label'], AdaBoostPrediction)
print("Ada Boost Accuracy:",AdaBoostAccuracy)
print("Confusion Matrix:")
print(confusion_matrix(test_news['Label'], AdaBoostPrediction))
```

```
Ada Boost Accuracy: 0.5633802816901409
Confusion Matrix:
[[125 115]
 [102 155]]
```

Fig 3.1 AdaBoost Classifier Result

After building the AdaBoost classifier model, the accuracy of the model is calculated. The performance of the AdaBoost classifier model is evaluated using the confusion matrix.

Table -1: Results

S.N o.	ML MODEL	ACCURACY
1	Decision Tree Algorithm	51.1%
2	SGD Algorithm	53.11%
3	Logistic Regression	53%
4	Random Forest Classifier	54.2%
5	KNN Classifier	55.2%
6	AdaBoost Classifier	56.3%

3. CONCLUSION

Using only textual datasets, the proposed system with various algorithms to predict stocks while comparing them with each other. Classification with the numerical data has been done using linear regression. In our project, Adaboost classifier gives higher accuracy of 56.3% compared to other algorithms

REFERENCES

- [1] https://webfocusinfocenter.informationbuilders.com/wfappent/TLS/TL_rstat/ source/LogisticRegression43.htm
- [2] <http://dsbyprateekg.blogspot.com/2017/09/machine-learning-decision-trees- and.html>
- [3] <https://www.geeksforgeeks.org/ml-stochastic-gradient-descent-sgd/>
- [4] <https://towardsdatascience.com/boosting-algorithm-adaboost-b6737a9ee60c>
- [5] <https://www.kaggle.com/aaron7sun/stocknews>
- [6] https://www.tutorialspoint.com/machine_learning_with_python/machine_learning_with_python_knn_algorithm_finding_nearest_neighbors.htm