Prediction in Market Volatility

A case study in predicting market volatility and building short-term trading strategies using data from Reddit's WallStreetBets.

Pavan kumar C N

CONTENT

- Project Approach
- What does the DATA tell us?
- Our **PREDICTION** models -**PERFORMANCE** evaluate
- **CONCLUSION & NEXT STEP**

Helps to make a prediction on stock prices and market volatility.

SCENARIO

The aim of this project is to use data from posts made made on the sub-reddit "Wallstreet-Bets" to make a prediction of given scenario.

Covers two datasets:

JSON file: - Contains comment of Reddit's post. - Performed

Sentiment

Help to predict if specific stocks rose or fell in the given time frame.

SCENARIO

Analysis.

Excel file:

 Trimmed this huge org. provided data as per the other similar file hosted on Kaggle.

Predict Market Volatility

, why?

How can predicting market volatility add values to business world. Current scenarios' relation between stock market and social media.

Sudden market volatility increment affects the investment so predicting Market volatility in advance can increase /lead us to profits in Stock market.

80% of investors today use it as their regular Workflow &

Approx. 30% obtain information

about the investment market through different

Social Media (it).

Economic Times
Research

Target Variable

- Created comparing today's close price and yesterday's closing price.
- Check how the sentiment analysis of comment made in the day affects the closing price.

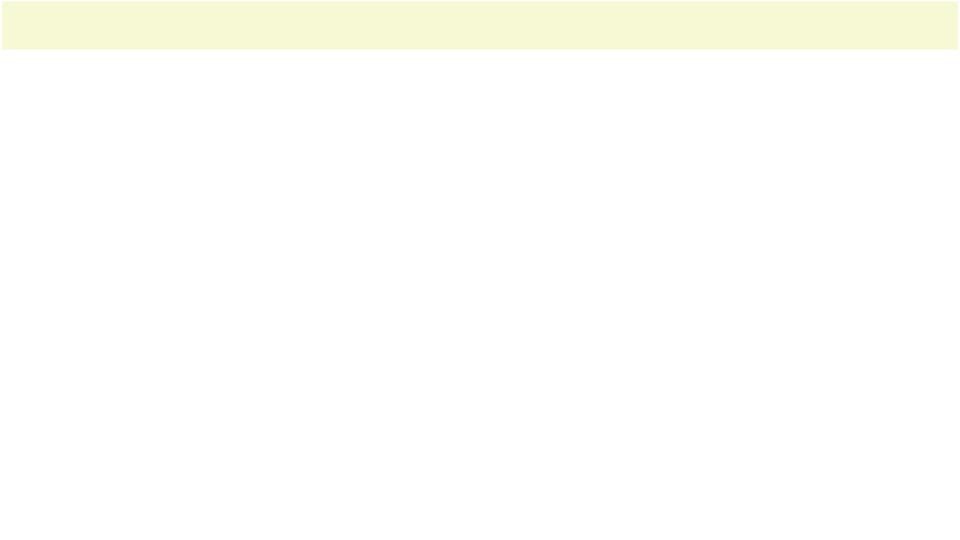
If we Predict the can future profit/
loss, we can AVOID

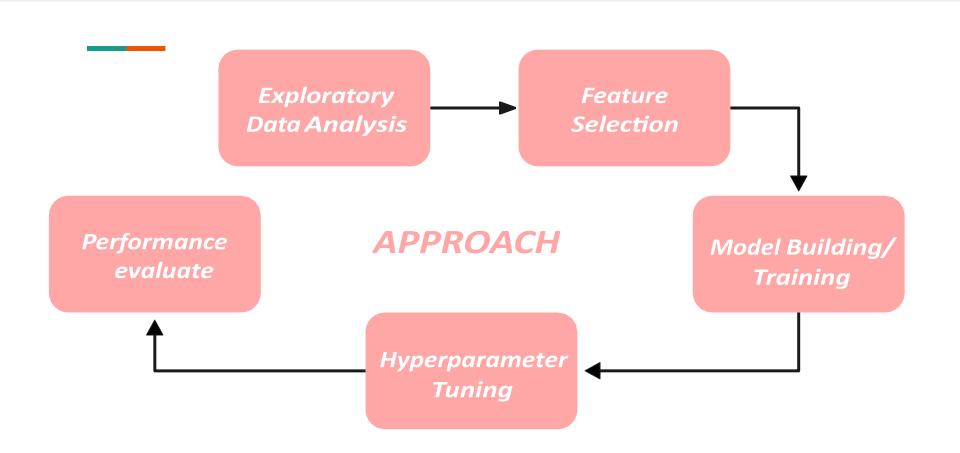
Why
Data
Science?

the market volatility and get maximum profit from current stocks.

APPROACH

How data science helps to predict the market volatility, and how we are going to do with it.



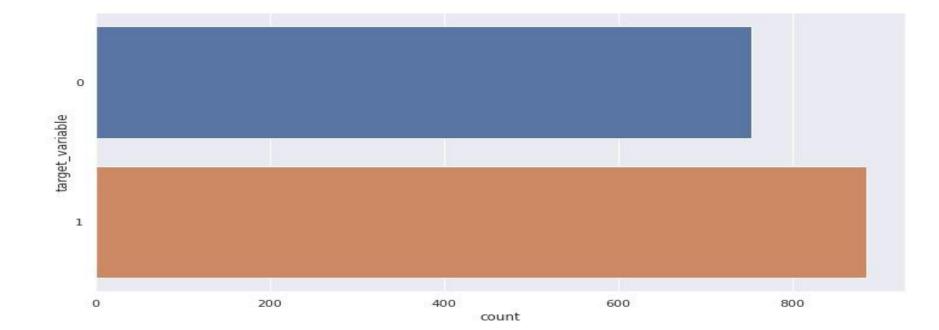


Data Analysis

What the data told us? Let go for an EDA on the data set.

Our Target Variable (P/L)

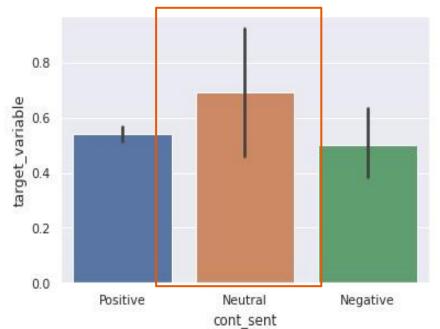
More positive response/profit in datasets.



Neutral >
Positive>
Negative
Responses
seems to affect
target variable.

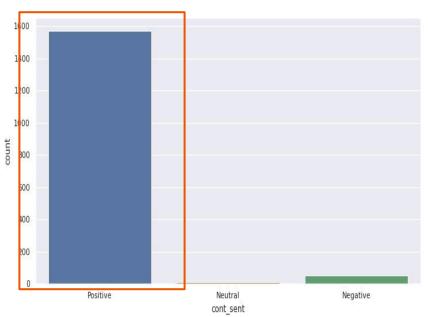


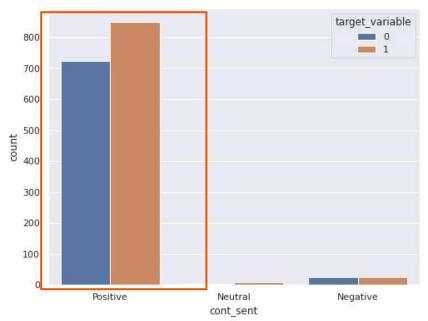
Bi-variate Plot (Polarity-Sentiment Analysis)



Count Plot (Univariate and Bivariate)

- Univariate Plot (Number of positive responses > Negative > Neutral)
- Positive response influences profit(1) in target variable/Closing Price.





Relation Among all the **Dependent** And Independent **Variables**

Heatmap

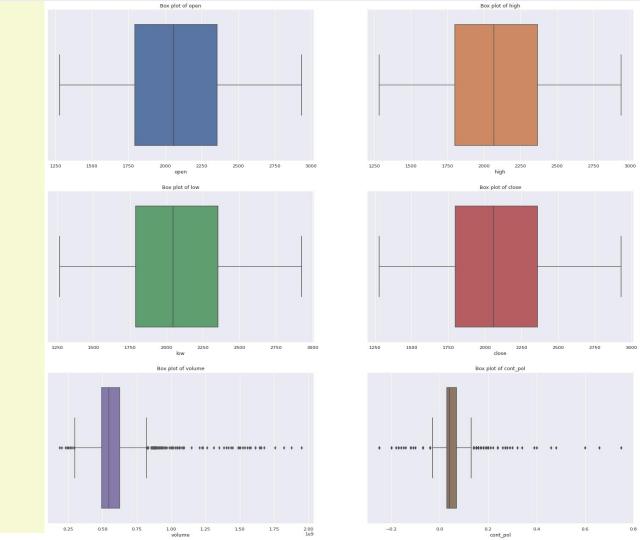
| | | | | | | | | | | | | | | | lia: | |
|--|---------|---------|----------|---------|---------|----------|----------|-----------|------------|----------|----------|----------|----------|-----------------|--------|---|
| 5 8 | 1 | 1 | 1 | 1 | -0.0067 | -0.46 | -0.27 | 0.96 | -0.0019 | 0.0068 | -0.22 | -0.13 | 0.26 | -0.015 | -1.00 | I |
| Sept. | 1 | | | 1 | -0.0033 | | -0.27 | 0.96 | -0.0033 | 0.0062 | -0.22 | -0.13 | 0.26 | -0.0045 | - 0.75 | I |
| o o | 1 | | | 1 | -0.016 | | -0.27 | 0.96 | 0.00043 | 0.0078 | -0.22 | -0.13 | 0.26 | -0.00069 | | I |
| | 1 | 1 | 1 | i | -0.011 | -0.47 | -0.27 | 0.96 | -0.0019 | 0.0058 | -0.22 | -0.13 | 0.26 | 0.0081 | - 0.50 | I |
| wolum e | -0.0067 | -0.0033 | -0.016 | -0.011 | 1 | -0.052 | -0.054 | 0.067 | -0.089 | 0.075 | -0.018 | -0.013 | 0.022 | -0.094 | | I |
| iii ta saa saa saa saa saa saa saa saa saa | | | | -0.47 | -0.052 | 1 | 0.19 | -0.51 | 0.051 | 0.0056 | 0.033 | 0.0029 | -0.031 | -0.013 | - 0.25 | I |
| 10 To | | -0.27 | -0.27 | -0.27 | -0.054 | 0.19 | 1 | -0.31 | 0.014 | -0.013 | -0.38 | -0.086 | 0.38 | 0.0076 | | I |
| re vear | 0.96 | 0.96 | 0.96 | 0.96 | 0.067 | -0.51 | -0.31 | 1 | -0.14 | -0.007 | -0.22 | -0.12 | 0.26 | -0.0048 | - 0.00 | I |
| though the state of the state o | -0.0019 | -0.0033 | 0.00043 | -0.0019 | -0.089 | 0.051 | 0.014 | -0.14 | 1 | -0.0071 | 0.014 | -0.016 | -0.0054 | -0.021 | | I |
| APD APP APP APP APP APP APP APP APP APP | | 0.0062 | 0.0078 | 0.0058 | 0.075 | 0.0056 | -0.013 | -0.007 | -0.0071 | 1 | -0.014 | 0.011 | 0.0079 | -0.045 | 0.2 | i |
| Day 1,500 Bay 1, | | -0.22 | -0.22 | -0.22 | -0.018 | 0.033 | -0.38 | -0.22 | 0.014 | -0.014 | 1 | -0.016 | -0.89 | -0.015 | 0.5 | |
| nau saut | -0.13 | -0.13 | -0.13 | -0.13 | -0.013 | 0.0029 | -0.086 | -0.12 | -0.016 | 0.011 | -0.016 | 1 | -0.44 | 0.027 | 0.3 | |
| sod Just | 0.26 | 0.26 | 0.26 | 0.26 | 0.022 | -0.031 | 0.38 | 0.26 | -0.0054 | 0.0079 | -0.89 | -0.44 | 1 | 0.00088 | 0.7 | 5 |
| variable | | -0.0045 | -0.00069 | | -0.094 | -0.013 | 0.0076 | | | -0.045 | -0.015 | 0.027 | 0.00088 | 1 | | |
| lander veri | uado | high | low | close | volume | cont_len | cont_pol | date_year | date_month | date_day | sent_neg | sent neu | sent_pos | target_variable | | |

Check the Outliers:

All the variables are outliers free other than:

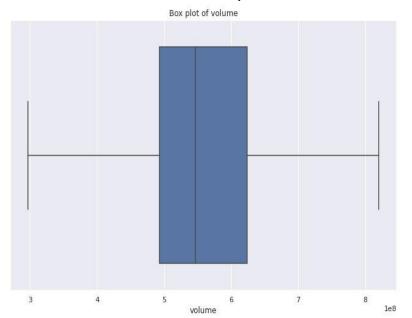
- Volume
- Content Polarity

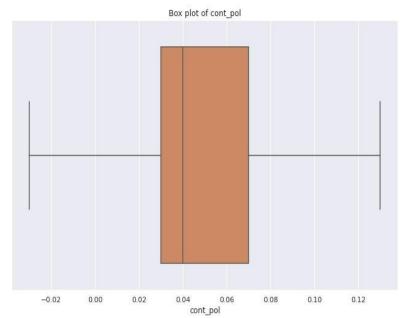
BoxPlot



Box Plot (After):

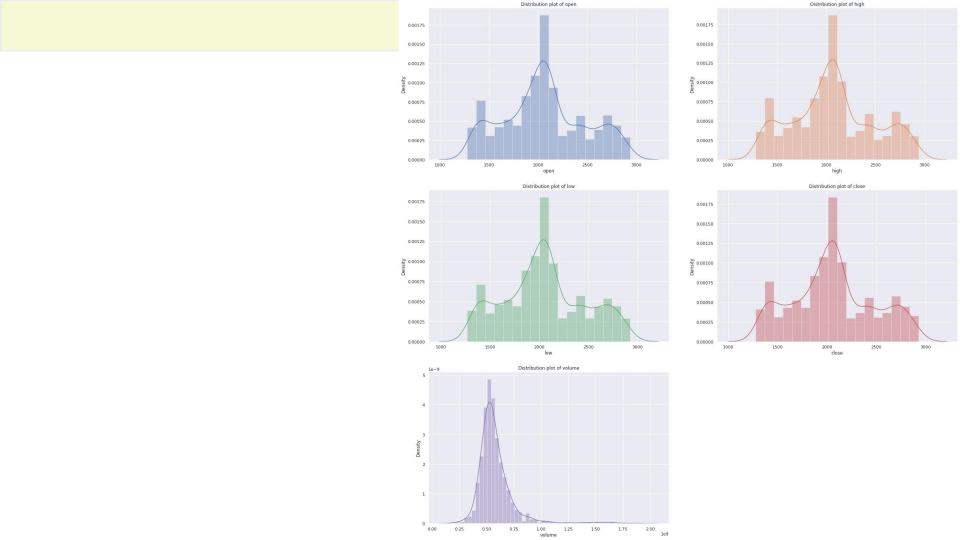
- IQR was performed where the outliers were treated via flooring and capping.
- Volume and Content Polarity columns are now outliers free.





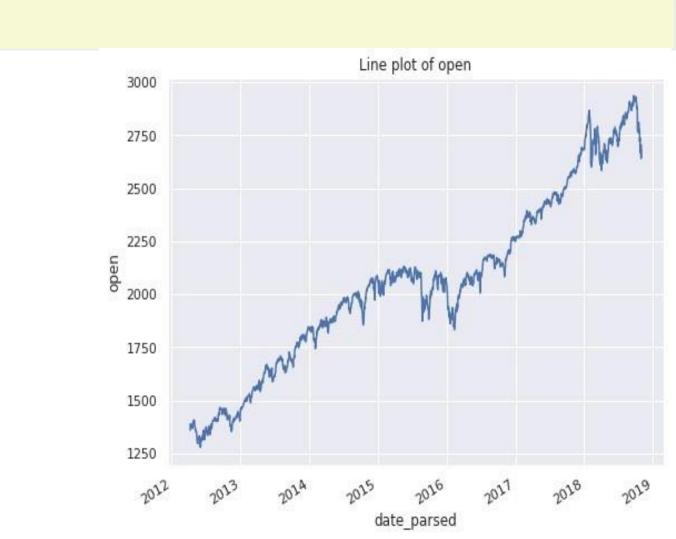
Data of all the variables are Normally Distributed.

Disribution Plot



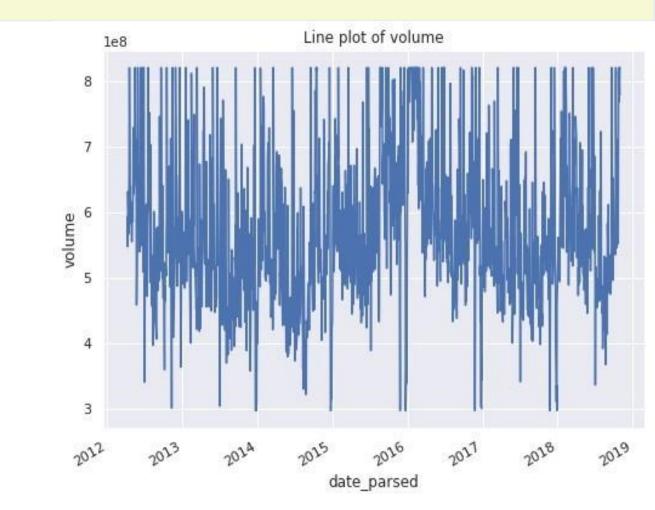
- It seems that all the other variables like: Close, High, Low has similar line plot other than Volume.
- All the variables value seems to increase as per the time.

Line Plot



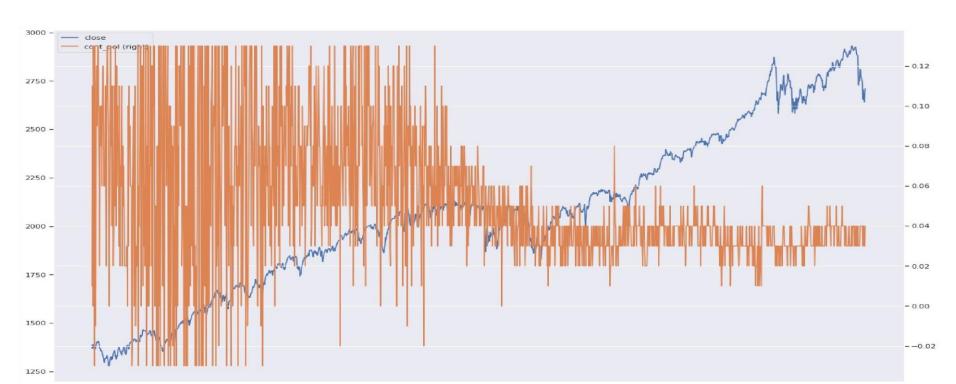
Volume seems to increase and decrease along with time.

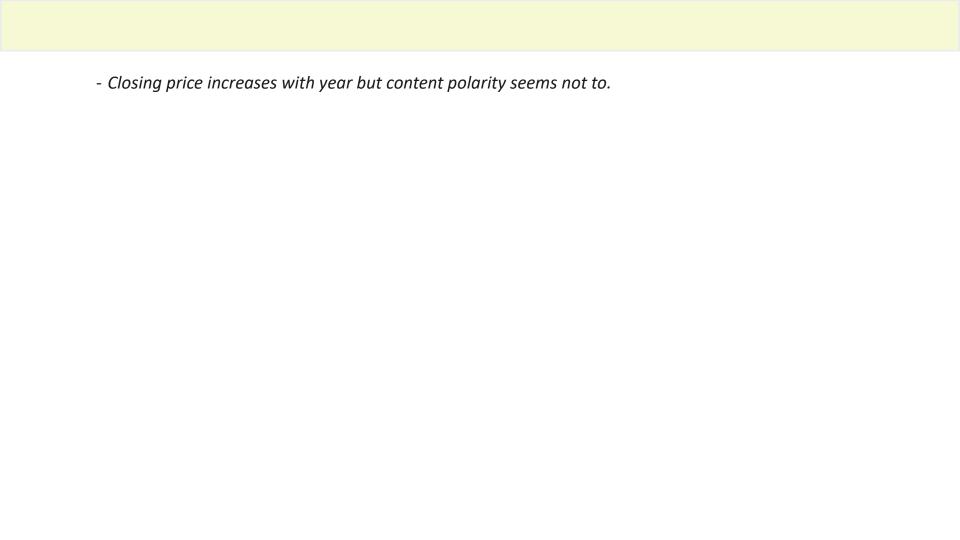
Line Plot



Box Plot (After):

- Seems like closing price and Content Polarity are correlated with one another.





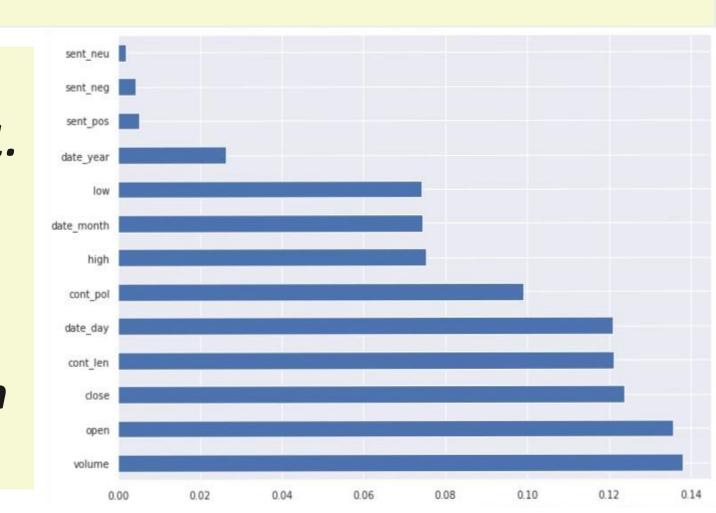
PREDICTION MODEL

Build a classification model to predict the market volatility.

Most Imp. Features: 1. Volume 2. Open

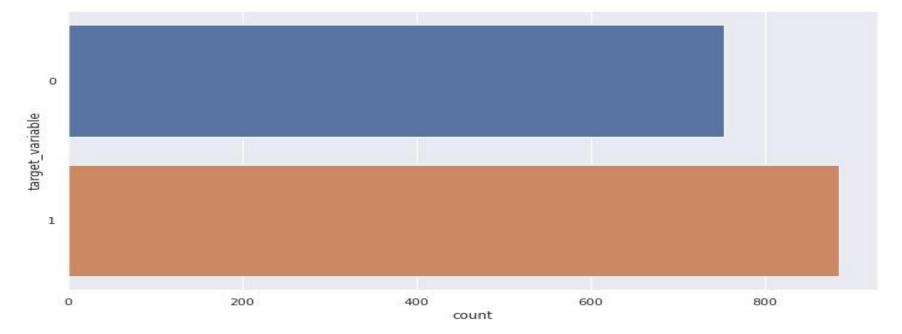


5. Date



TARGET VARIABLE

- Though positive target variable is quite more in comparison to negative target variable we cannot say dataset is imbalanced because data the difference is not so huge.



MODEL Building/Training

Logistic Regression was selected for a model.

```
log_reg.fit(x_train, y_train)

y_pred = log_reg.predict(X_test)
print(classification_report(y_test, y_pred))

acc_score = accuracy_score(y_test,y_pred)
acc_score_per = acc_score * 100
print('The accuracy score is', acc_score, '/', acc_score_per, '%'.)

MODEL BUILDING — Logistic Regression
```

from sklearn.learn model import LogisticRegression

- Classification Report and Accuracy score of our model (Before Hyperparameter Tuning)

| port | suppo | T1-score | recall | precision | |
|------|-------|----------|--------|-----------|-------------|
| 161 | 1 | 0.57 | 0.43 | 0.83 | 0 |
| 167 | 1 | 0.74 | 0.92 | 0.62 | 1 |
| 328 | 3 | 0.68 | | | accuracy |
| 328 | 3 | 0.65 | 0.67 | 0.73 | macro avg |
| 328 | 3 | 0.66 | 0.68 | 0.73 | eighted avg |

The accuracy score is 0.676829268292683 / 67.6829268292683 %.

PERFORMANCE EVALUATION

Hyperparameter Tuning/ Evaluation metrics to increase the accuracy of the model.

Confusion Matrix

True Negative: 69 (Predicted Loss as Loss)

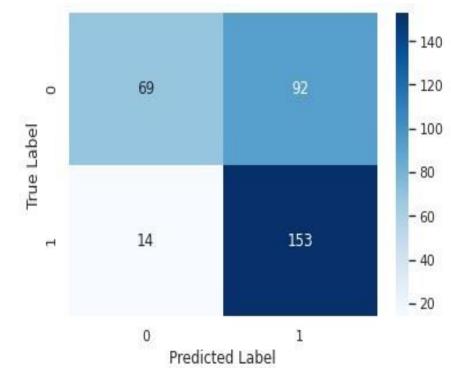
False Positive: 92 (Predicted Loss as Profit)

False Negative: 14 (Predicted Profit as Loss)

True Positive: 153 (Predicted Profit as Profit)

Summary





AOC-ROC Curve

Before Hyperparameter Tuning

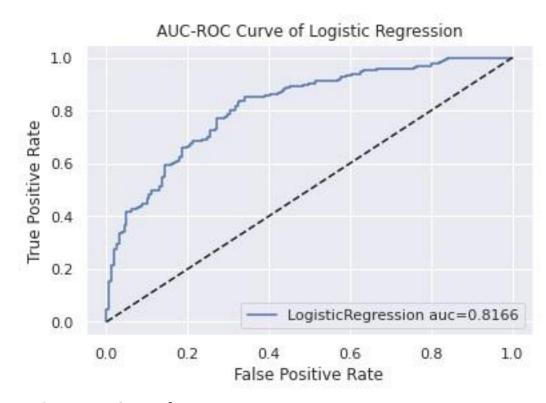
ROC Score

60.8166400119016626 / 81.66400119016626 %

Graph

The left corner of model is quite near to top-left corner but not exactly so the roc curve of is average.

Summary



Hyperparameter Tuning (GridSearchCV)

```
from sklearn.model selection import GridSearchCV
penalty=['11', '12', 'elasticnet']
solver=['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']
max iter=[100, 200, 300, 350]
random grid={ 'penalty':penalty,
'solver':solver,
             'max iter':max iter,
log reg grid search= GridSearchCV(estimator=log reg,
param grid=random grid, cv=20, n jobs=-1, verbose=2)
```

Confusion Matrix

True Negative: 148

(Predicted Loss as Loss)

False Positive: 13

(Predicted Loss as Profit)

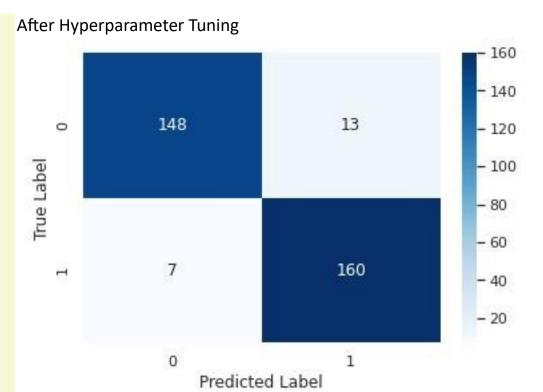
False Negative: 7

(Predicted Profit as Loss)

True Positive: 160

(Predicted Profit as Profit)

Summary



AOC-ROC Curve

ROC Score

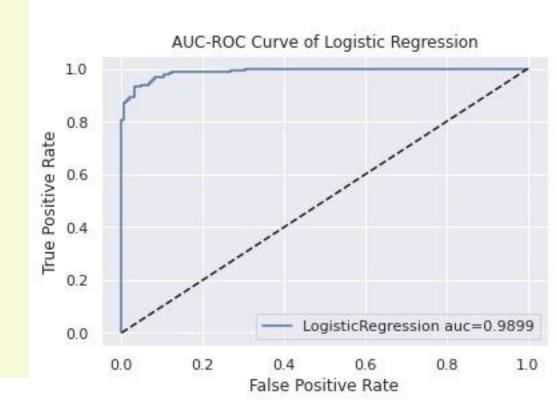
0.98988358686354 / 98.988358686354 %

Graph

The left corner of model is so close to top-left corner hence model is good.

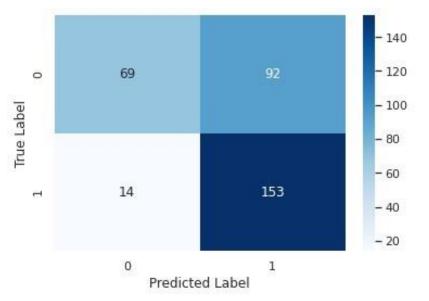
Summary

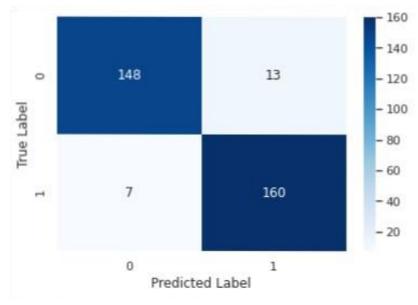
After Hyperparameter Tuning



MODELS PERFORMANCES

- Classification Report and Accuracy score of our model (After Hyperparameter Tuning) - Increment in True Positive/ False Positive as expected..





MODEL BUILDING – Logistic Regression

- Classification Report and Accuracy score of our model (After Hyperparameter Tuning)

| р | recision | recall | f1-score | support | |
|------|----------|--------|----------|---------|--|
| Θ | 0.95 | 0.92 | 0.94 | 161 | |
| 1 | 0.92 | 0.96 | 0.94 | 167 | |
| racy | | | 0.94 | 328 | |
| avg | 0.94 | 0.94 | 0.94 | 328 | |
| avg | 0.94 | 0.94 | 0.94 | 328 | |

Model Deployment

Flask along with HTML/CSS was used to deploy in local server. Later deployed using Heroku.

MODEL DEPLOYMENT

Tools:

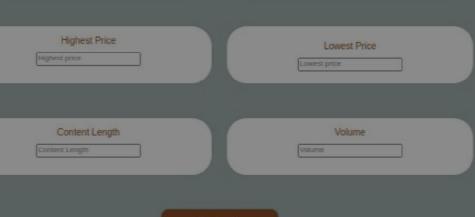
Flask, HTML, CSS, Heroku

Input

Date, Sentiment Analysis, Open, Close, Higher Price, Lower Price, Content Length, Close

Summary

Date Sentiment Analysis (Polarity) Positive mm/dd/yyyy. --- --~ Open Close Highest Price Lowest Price



MODEL DEPLOYMENT

Tools:

Flask, HTML, CSS, Heroku

Output

Gives the predicted output from the trained model in the form of Profit/Loss.



Result:



CONCLUSION

MODEL CONCLUSION:

CONDITIONS which have the following characteristics,

- Having HIGH opening price itself;
- High Volume;
- Positive Sentiment Analysis;
- Lengthy/Informative Detailed comments; are likely to

lead us to Profit.

WHAT CAN WE DO

General

- 1. Publish more *Positive Contents*;
- 2. Promote more **detailed and informative** contents
- 3. Reduce/Remove the **negative contents** from Social Media asap if found.

LIMITATION & NEXT STEP

HOW TO IMPROVE

1. Only applied Logistic Regression:

Apply and compare other tuned performance.

2. Used only Reddit's API:

Collect API from as much as resources possible.

END