

- I have selected **Naive Bayes** as classification approach for supervised learning but Multinomial bayes classifier is used for model training. It is expected to give better performance over regression logistic based due to following reasons:
 - Bayes theorem finds the probability of event based on event already occurred
 - It check the probability of even after the evidence is seen.
 - Naive assumption are taken in the naive bayes classifier to split the evidence into distint parts
 - Multinomial Naive Bayes: Feature vectors represent the frequencies with which certain events have been generated by a multinomial distribution.

Resource Reference:

1. <https://www.geeksforgeeks.org/naive-bayes-classifiers/?ref=gcse>

```
In [93]: # Importing Multinomial naive bayes from sklearn
from sklearn.naive_bayes import MultinomialNB
# Creating a classifier
nb = MultinomialNB()
# Training a model
nb.fit(X_tf_train, y_train) # Term frequency is provided for training set

Out[93]: MultinomialNB()

In [94]: # Testing model by providing testing dataset
y_pred = nb.predict(X_tf_test)

In [95]: import collections
# checking count of Sentiments in model predictions
collections.Counter(y_pred)

Out[95]: Counter({'positive': 224652, 'negative': 230307})

In [96]: collections.Counter(y_test)

Out[96]: Counter({'positive': 225772, 'negative': 229187})
```

- Accuracy of Naive bayes

```
In [98]: # Checking the Accuracy of Model
from sklearn import metrics
print('Naive Bayes accuracy (s):', metrics.accuracy_score(y_test, y_pred)*100)

Naive Bayes accuracy (s): 77.3764229304179

We can notice that naive bayes accuracy (i.e. 77.376) is superior in comparison to Logistic regression (68.18).
```

- Confusion Matrix of Naive bayes



Accuracy of both models are further analysed in the section below using "Precision, Recall, F1-Measure and ROC" matrices.

Section 3. Conclusion

3.1 | Precision, Recall & F1-Measure

- **Accuracy** checks the correct prediction made by model over the total observation
- **Precision** checks the count of correct prediction made by the model
- **Recall** is opposit to the precision. It check how often correct prediction made by the model when result is actually correct
- **F-measure** is harmoc mean of both precision and recall

Formulas

- **Accuracy**: Number of correct predictions / total number of predictions ==> (TP + TN) / (TP + TN) + FP + FN)
- **Precision**: (TP) / (TP + FP)
- **Recall**: (TP) / (TP + FN)
- **F1 Measure**: 2 * (precision * recall) / (precision + recall)

Resource Reference:

1. <https://towardsdatascience.com/understanding-confusion-matrix-99a424d4d62>
2. <https://vitalflux.com/accuracy-precision-recall-f1-score-python-example/>

Since, we have already calculated accuracy of both classifier above, therefore, these are not reproduced or recalculated in this section.

- **Baseline performance** - Precision, Recall and F1 Measure of linear regression

```
Precession Score

In [101]: #Checking the precession score of logistic regression
print('Logistic Regression Precision score: %.3f' % metrics.precision_score(y_test, lg_y_pred, pos_label='positive'))
Logistic Regression Precision score: 0.684

Recall Score

In [102]: #Checking the recall score of logistic regression
print('Logistic Regression Recall score: %.3f' % metrics.recall_score(y_test,lg_y_pred, pos_label='positive'))
Logistic Regression Recall score: 0.689

F1 Measure

In [103]: #Checking F1 measure of Logistic regression
print('Logistic Regression F1 measure : %.3f' % metrics.f1_score(y_test, lg_y_pred, pos_label='positive'))
Logistic Regression F1 measure : 0.687

• Classification approach - Precision, Recall and F1 Measure of Naive Bayes

Precession Score

In [104]: #Checking the precession score of Naive Bayes
print('Naive Bayes Precision score: %.3f' % metrics.precision_score(y_test, y_pred, pos_label='positive'))
Naive Bayes Precision score: 0.773

Recall Score

In [105]: #Checking the recall score of Naive Bayes
print('Naive Bayes Recall score: %.3f' % metrics.recall_score(y_test,y_pred, pos_label='positive'))
Naive Bayes Recall score: 0.770

F1 Measure

In [106]: #Checking F1 measure of Naive Bayes
print('Naive Bayes F1 measure : %.3f' % metrics.f1_score(y_test, y_pred, pos_label='positive'))
Naive Bayes F1 measure : 0.771

Comparison between Logical Regression and Multinomial Naive Bayes model
```

	Logical Regression	Multinomial Naive Bayes
Accuracy	0.688	0.774
Precision	0.684	0.773
Recall	0.689	0.770
F-1 Measure	0.687	0.771

3.2 | ROC Curve

- **ROC curve (receiver operating characteristic curve)** it is a graph to show the performance of Classifier at classification thresholds. It is based on two parameters:
 - **True Positive Rate**: synonym for recall ==> TPR = (TP) / (TP + FN)
 - **False Positive Rate**: FPR = (FP) / (FP + TN)

Resource Reference:

1. <https://developers.google.com/machine-learning/crash-course/classification/roc-and-auc>



ROC curve showing comparison between Multinomial Naive Bayes and Regression logistic classifier. We can notice that AUC (Area under the ROC Curve) shows that Naive Bayes was efficient since beginning and it has more probability of ranking random positive than a random negative. AUC rate of Naive Bayes is 0.86 in comparison to Logistic Regression (which is referred as "Pipeline" in chart above)

3.3 | Summary and Conclusions

Conclusion

The purpose of this project is to prepare a text classifier, which could be beneficial for the investor or financial institution in assessing the public financial behaviour in order to align their investment decisions accordingly. Here is a quick recap of the approach successfully followed:

- Arrange a twitter dataset;
- Conduct pre-processing of the dataset;
- Start with quick and simple model;
- Get the classification prediction from base model;
- Prepare another classifier using a different approach to match the performance with the base model.

In terms of comparison between baseline and classification approach, Multinomial Naive Bayes has clearly outperformed Logical regression in categorizing the tweets as True Positive or True Negative. The Accuracy of Naive Bayes is 12.5% more in comparison to Logical Regression. Moreover, Naive Bayes has also posted better scores also for precision, recall and F-1 Measure.

Furthermore, ROC-AUC also shows Naive Bayes with an efficiency ratio of 0.86, and it has better categorization since start. Naive Bayes performed well in comparison, but the accuracy ratio of 77.4% is less than other models available online. It could be because of following reasons:

- The datasets used to train other models are quite small, and text classifiers such as Naive Bayes perform quite well on small datasets.
- As we have seen in the features extraction and presentation part, the same terms are used very often among positive and negative sentiments.

There is no standard benchmark available to judge any classifier, but there is always room for improvement on any implementation.

Further Improvement

- The performance and accuracy could be measured by using other classifiers of the Naive Bayes family, such as Gaussian Naive Bayes and Bernoulli Naive Bayes. These may model them more accurately and perform better as a result.
- Secondly, There is also the possibility of getting better results by changing re-processing methodology that converts text to numbers.

~ 313 Word Count

General Resource Reference:

Plotting Guideline:

- <https://towardsdatascience.com/different-bar-charts-in-python-6d984b9c6b17>
- <https://blog.inightsdatascience.com/how-to-solve-90-of-nlp-problems-a-step-by-step-guide-fda605278e4e>

Markdown Guideline:

- <https://www.markdownguide.org/basic-syntax/>
- <https://www.markdownguide.org/hacks/#color>

Literature:

- Speculator and Influencer Evaluation in Stock Market by Using Social Media [2020 IEEE publication - International Conference on Big Data (Big Data)]
- Sentiment Analysis in Financial Markets [2014 IEEE publication - University of Siegen Institute of Knowledge Based Systems Germany]