

# The Occidental Computer Science Comprehensive Project: Tutorial Report

Odelia Putterman

putterman@oxy.edu

Occidental College

## Abstract

This report documents the tutorial completed as part of my Occidental College Computer Science Comprehensive Project: *Predicting Cryptocurrency Prices for Stock Trading Using Machine Learning*. This report has four components: methods, evaluation, discussion, and software documentation. For each component, we include a section. This report walks us through the tutorial followed, Stock Market Sentiment Analysis Using Python & Machine Learning (Youtube tutorial), reviewing the key concepts and learned information relevant to the comprehensive project.

## 1 Methods

This section is a walk-through of the methods implord in this tutorial in a numbered list. This tutorial was done in Google Colab. We detail the steps below.

1. First, we installed *vaderSentiment* with pip to get access to sentiment analysis tools. These will be used later on in the tutorial.

```
1 pip install vaderSentiment
```

Listing 1. Install vaderSentiment

2. Next, we load all the necessary libraries for this tutorial. These are:

- pandas;
- numpy;
- textblob.TextBlob;
- re;
- vaderSentiment.SentimentIntensityAnalyzer;
- sklearn.model\_selection.train\_test\_split;
- sklearn.metrics.accuracy\_score;
- sklearn.metrics.classification\_report; and
- sklearn.discriminant\_analysis.LinearDiscriminantAnalysis.

```
1 #Import the libraries
2 import pandas as pd
3 import numpy as np
4 from textblob import TextBlob
5 import re
6 from vaderSentiment.vaderSentiment import
  SentimentIntensityAnalyzer
```

```
7 from sklearn.model_selection import
  train_test_split
8 from sklearn.metrics import accuracy_score,
  classification_report
9 from sklearn.discriminant_analysis import
  LinearDiscriminantAnalysis
```

Listing 2. Load Libraries

3. Next, we load in all the data for this tutorial using **google.colab.files.upload()**. There were two input data sets, which are the *Dow Jones Industrial Average News.csv* and the *Dow Jones Industrial Average Stock.csv*, all originating from data from the Dow Jones and downloaded from Kaggle at this link: <https://www.kaggle.com/datasets/aaron7sun/stocknews>.

```
1 #Load the data
2 from google.colab import files
3 files.upload()
```

Listing 3. Upload data sets

4. Next, we merged these two data sets into a pandas data frame called *merge* on the 'Date' field to consolidate our data. Now, this *merge* data frame had the original news headers, label, date, open, high, low, close, volume, and adjacent close columns.

```
1 #Store the data into variables
2 df1 = pd.read_csv('
  Dow_Jones_Industrial_Average_News.csv')
3 df2 = pd.read_csv('
  Dow_Jones_Industrial_Average_Stock.csv')
4
5 #Merge the data set on the date field
6 merge = df1.merge(df2, how='inner', on='Date'
  )
```

Listing 4. Store and merge data

5. To prepare the data for sentiment analysis, we grabbed the *news* data set inputs, consolidated them into one joint string for each date, and cleaned the data by removing certain unnecessary string patterns, storing this processed data in a new list and adding it to our consolidated csv under a new column called 'Combined\_News'. Now, the merged data frame had the original news headers, label, date, open, high, low, close, volume, and adjacent close columns in addition to newly added 'Combined\_News' column.

```

1 #Combine the top news headlines
2 headlines = []
3
4 for row in range(0, len(merge.index)):
5     headlines.append(' '.join( str(x) for x in
6                               merge.iloc[row, 2:27]))
7
8 #Clean the data
9 clean_headlines = []
10
11 for i in range(0, len(headlines)):
12     clean_headlines.append(re.sub("b[()]", '',
13                                   headlines[i])) # remove b'
14     clean_headlines[i] = re.sub('b[()]",', '',
15                                   clean_headlines[i]) # remove b"
16     clean_headlines[i] = re.sub("\\", '',
17                                   clean_headlines[i]) #remove \
18
19 #Add the clean headlines to the merge data
20 set
21 merge['Combined_News'] = clean_headlines

```

Listing 5. Pre-process data

- Next, we created functions to get the subjectivity and polarity of text using `TextBlob(text).sentiment.subjectivity` and `TextBlob(text).sentiment.polarity`, where `text` is replaced with the actual text to analyze.

```

1 #Create a function to get the subjectivity
2 def getSubjectivity(text):
3     return TextBlob(text).sentiment.
4         subjectivity
5
6 #Create a function to get the polarity
7 def getPolarity(text):
8     return TextBlob(text).sentiment.polarity

```

Listing 6. Subjectivity and polarity functions

- We used these newly created functions to create a subjectivity and polarity column, called 'Subjectivity' and 'Polarity' respectively, with inputs as the subjectivity and polarity outputs from these functions on each line of 'Combined\_News' (the cleaned news string).

```

1 #Create two new columns 'Subjectivity' and '
2   Polarity'
3 merge['Subjectivity'] = merge['Combined_News']
4   .apply(getSubjectivity)
5 merge['Polarity'] = merge['Combined_News']
6   .apply(getPolarity)

```

Listing 7. Get subjectivity and polarity

- Next, we made a function to get the sentiment scores called `getSIA`. Using this function, we got the sentiment scores for each date, creating a new column in the combined csv for each of: 'combined', 'negative', 'neutral', and 'positive', where each is a value between -1 and 1.

```

1 #Create a function to get the sentiment
2   scores
3 def getSIA(text):
4     sia = SentimentIntensityAnalyzer()
5     sentiment = sia.polarity_scores(text)

```

```

5     return sentiment
6
7 #Get the sentiment scores for each day
8 compound = []
9 neg = []
10 pos = []
11 neu = []
12 SIA = 0
13
14 for i in range(0, len(merge['Combined_News']))
15     ):
16     SIA = getSIA(merge['Combined_News'][i])
17     compound.append(SIA['compound'])
18     neg.append(SIA['neg'])
19     neu.append(SIA['neu'])
20     pos.append(SIA['pos'])
21
22 #Store the sentiment scores in the merge data
23 set
24 merge['Compound'] = compound
25 merge['Negative'] = neg
26 merge['Neutral'] = neu
27 merge['Positive'] = pos

```

Listing 8. Add sentiment scores

- We cleaned our combined csv to keep only the necessary columns: 'Open', 'High', 'Low', 'Volume', 'Subjectivity', 'Polarity', 'Compound', 'Negative', 'Neutral', 'Positive', and 'Label'.

```

1 #Create a list of columns to keep
2 keep_columns = ['Open', 'High', 'Low', '
3   Volume', 'Subjectivity', 'Polarity', '
4   Compound', 'Negative', 'Neutral', '
5   Positive', 'Label']
6
7 df = merge[keep_columns]

```

Listing 9. Keep only necessary columns

- We split this cleaned combined csv into two data frames, one for the 'input' data (everything but the 'Label' column) and one for the target data (the 'Label' columns).

```

1 #Create the feature data set
2 X = df
3 X = np.array(X.drop(['Label'], 1))
4
5 #Create the target data set
6 y = np.array(df['Label'])

```

Listing 10. Split into input and target data

- Using `train_test_split`, we split the data into 80 percent and 20 percent to train and test the data.

```

1 #Split the data into 80% training and 20%
2   testing data sets
3 x_train, x_test, y_train, y_test =
4   train_test_split(X, y, test_size = 0.2,
5   random_state = 0)

```

Listing 11. Split data into training and test

- Using the train data, we trained a `LinearDiscriminantAnalysis()` model.

```

1 #Create and train the model

```

	precision	recall	f1-score	support
0	0.86	0.79	0.83	193
1	0.82	0.88	0.85	205
accuracy			0.84	398
macro avg	0.84	0.84	0.84	398
weighted avg	0.84	0.84	0.84	398

Figure 1. Evaluation metrics for tutorial results.

```
2 model = LinearDiscriminantAnalysis().fit(
    x_train, y_train)
```

Listing 12. Train model

13. And, finally, we got the model predictions for the test data and compared between these predictions and the labeled data to get the precision, recall, f1-score, and support.

```
1 #Show the models predictions
2 predictions = model.predict(x_test)
3
4 #Show the model metrics
5 print(classification_report(y_test,
    predictions))
```

Listing 13. Get predictions and evaluate

## 2 Evaluation

This tutorial evaluated its outputted results by comparing the predicted stock prices produced from the trained model with the actual stock prices from that same period. This found to produce a precision rate over 80%. Other measures besides precision were accounted for as well, including recall, f1-score, and support, all showing very promising results.

To evaluate our results, we split the data into 80% training data and 20% testing data using **train\_test\_split** from **sklearn.model\_selection**. After training the model using **LinearDiscriminantAnalysis()** on the training data, we used the model to make predictions for this test data. Next, we pulled the actual price changes from this time (on the 0/1 scale) and used **classification\_report** from **sklearn.metrics** to get the precision, recall, f1-score, support, accuracy, macro average, and weighted average, which are all metrics we can use to evaluate the success of our model which were, as mentioned, highly promising.

## 3 Discussion

Considering the relatively small size of the initial inputted data sets (less than 5.5 MB combined), I am very impressed with the high accuracy rates. Moreover, the tutorial was only trained on 80% of this input data, meaning even less data was used in getting these impressive results.

One the one hand, this tutorial leads me to believe that with the right input data, you do not need an enormous data set, but on the other hand, it suggests that if we collect even more relevant data, we can get even more outstanding results.

All that said, I have a few critiques for this tutorial's input formatting and outputted results. The result outputs simply gave a '0' or '1', where zero signified decreasing stock price and one stood for increasing stock price. This zero or one system is not very insightful, as we do not know *how much* the stock is increasing or decreasing in price. If we could replace this system with a percentage of growth or decay, we could get more insightful data for the accuracy rates of our model but also for what to do with the results. As such, I would opt to switch the input stock price growth column with a numeric percentage in place of its current zero or one ranking. In doing so, our outputs would vary dramatically in projected percentage increase or decrease, allowing us to make more educated decisions and better evaluate the model's true accuracy.

Another critique I have is for the input data processing. In pre-processing the inputted file *DowJonesIndustrialAverageNews.csv*, we processed the file to remove unnecessary, often repeated characters. Though this worked in this context because we knew the exact format of the input data, I would suggest using a preset list of words or character symbols to filter out on instead so we can account for unforeseen but potentially problematic tokens.

This tutorial walked me through how to predict stock prices (whether they will increase or decrease) with sentiment analysis based on top news headlines. The whole process was so seamless and straight-forward it boosted my confidence in my ability to carry-out this project tremendously. While it showed the ease with which sentiment analysis and model building can be done using pre-built libraries, it did not cover the complexity involved with sourcing data, as these data sets were simply provided in a form which was easy to work with. That said, data processing was covered, which is essential to the success of the model.

## 4 Software Documentation

See attached **stock\_market\_tutorial.pdf** file for the code from this tutorial.