

Reddit News Headlines and Stock Market Performance

In the current financial market, predicting daily stock prices using daily news as input is clearly a very complex task. Moreover, the market is very volatile because it is a result of multiple factors that change continuously.

Our research questions are centered around using reddit data to predict the Dow Jones stock average in two different ways, using text word count and using text sentiment of reddit headlines.

Our research questions are:

- 1. Can the text content of reddit daily news headlines be used to predict past stock prices?**
- 2. Can the sentiment of reddit daily news headlines be used to predict past stock prices?**

Load and Clean Data

Begin by loading and formatting the data.

Dataset includes the top 25 most viewed news stories from reddit for a day and the Dow Jones stock performance over this time.

In [172]:

```

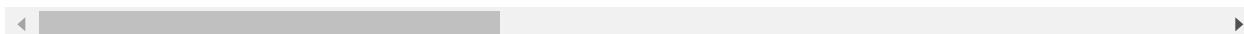
1 #
2 import pandas as pd
3
4 com_news_df = pd.read_csv('data/Combined_News_DJIA.csv')
5 com_news_df.head()

```

Out[172]:

	Date	Label	Top1	Top2	Top3	Top4	Top5	Top6	
0	2008-08-08	0	b"Georgia 'downs two Russian warplanes' as cou...	b'BREAKING: Musharraf to be impeached.'	b"Russia Today: Columns of troops roll into So...	b"Russian tanks are moving towards the capital...	b"Afghan children raped with 'impunity,' U.N. ...	b"150 Russian tanks have entered South Ossetia...	b"Br (i (
1	2008-08-11	1	b"Why wont America and Nato help us? If they w...	b'Bush puts foot down on Georgian conflict'	b"Jewish Georgian minister: Thanks to Israeli ...	b'Georgian army flees in disarray as Russians ...	b"Olympic opening ceremony fireworks 'faked'"	b'What were the Mossad with fraudulent New Zea...	b ang sale
2	2008-08-12	0	b'Remember that adorable 9-year-old who sang a...	b"Russia 'ends Georgia operation'"	b""If we had no sexual harassment we would hav...	b"Al-Qa'eda is losing support in Iraq because ...	b'Ceasefire in Georgia: Putin Outmaneuvers the...	b'Why Microsoft and Intel tried to kill the XO...	b' The G V th
3	2008-08-13	0	b' U.S. refuses Israel weapons to attack Iran:...	b"When the president ordered to attack Tskhinv...	b' Israel clears troops who killed Reuters cam...	b'Britain\'s policy of being tough on drugs is...	b'Body of 14 year old found in trunk; Latest (...	b'China has moved 10 *million* quake survivors...	ann Of Ge In
4	2008-08-14	1	b'All the experts admit that we should legalis...	b'War in South Osetia - 89 pictures made by a ...	b'Swedish wrestler Ara Abrahamian throws away ...	b'Russia exaggerated the death toll in South O...	b'Missile That Killed 9 Inside Pakistan May Ha...	b"Rushdie Condemns Random House's Refusal to P...	b ε (

5 rows × 27 columns



In [174]:

```

1 #DJIA_data = pd.read_csv('C:/Users/Yonas/OneDrive/Desktop/IMT 575/Group Proj
2 DJIA_data = pd.read_csv('data/upload_DJIA_table.csv')
3
4 #DJIA_data.info()

```

In [175]:

```

1 #Look at the shape of the data
2 print(com_news_df.shape, DJIA_data.shape )

```

(1989, 27) (1989, 7)

Merge the Dow Jones and the top news from reddit together. Also create a copy of the DF in order to use for sentiment analysis later.

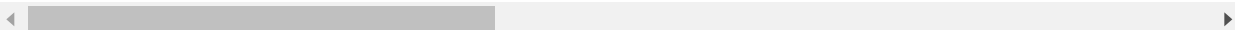
```
In [176]: 1 merged_data = com_news_df.merge(DJIA_data, how = 'inner', on = 'Date', left_
2
3 merged_data_sent = merged_data.copy())
```

```
In [177]: 1 merged_data.head(3)
```

Out[177]:

	Date	Label	Top1	Top2	Top3	Top4	Top5	Top6	
1988	2008-08-08	0	b"Georgia 'downs two Russian warplanes' as cou...	b'BREAKING: Musharraf to be impeached.'	b"Russia Today: Columns of troops roll into So...	b"Russian tanks are moving towards the capital...	b"Afghan children raped with 'impunity,' U.N. ...	b"150 Russian tanks have entered South Ossetia...	b"E
1987	2008-08-11	1	b"Why wont America and Nato help us? If they w...	b"Bush puts foot down on Georgian conflict'	b"Jewish Georgian minister: Thanks to Israeli ...	b"Georgian army flees in disarray as Russians ...	b"Olympic opening ceremony fireworks 'faked'"	b"What were the Mossad with fraudulent New Zea...	an
1986	2008-08-12	0	b'Remember that adorable 9-year-old who sang a...	b"Russia 'ends Georgia operation'"	b""If we had no sexual harassment we would hav...	b"Al-Qa'eda is losing support in Iraq because ...	b"Ceasefire in Georgia: Putin Outmaneuvers the...	b"Why Microsoft and Intel tried to kill the XO...	b

3 rows × 33 columns



Merge all 25 top headlines together

```
In [178]: 1 # to combine the first 25 coulums of the top news in
2
3 top25_headlines = []
4
5 for row in range(0, len(merged_data.index)):
6     top25_headlines.append(' '.join( str(x) for x in merged_data.iloc[row, 2
```

Optionally delete other extra columns of data [2:27]

```
In [181]: 1
2 #merged_data.drop(merged_data.iloc[:, 3:28], inplace = True, axis = 1)
3 #merged_data
```

Sort the columns by date

```
In [218]: 1 merged_data = merged_data.sort_values(by="Date")
2 #merged_data.head(2)
```

Hide ipython warning messages

```
In [183]: 1 %%javascript
2 (function(on) {
3   const e=$( "<a>Setup failed</a>" );
4   const ns="js_jupyter_suppress_warnings";
5   var cssrules=$( "#"+ns);
6   if(!cssrules.length) cssrules = $("<style id='"+ns+"' type='text/css'>div.ou
7   e.click(function() {
8     var s='Showing';
9     cssrules.empty()
10    if(on) {
11      s='Hiding';
12      cssrules.append("div.output_stderr, div[data-mime-type*='.stderr'] {
13    }
14    e.text(s+' warnings (click to toggle)');
15    on=!on;
16  }).click();
17  $(element).append(e);
18 })(true);
```

Hiding warnings (click to toggle)

Preprocessing the code

1 Next clean the data to make the machine learning models more effective.

Use regex to remove stopwords and to lemmatize all of the top 25 headlines.

```
In [184]: 1 import regex as re
2 from nltk.corpus import stopwords
3 from nltk.stem import WordNetLemmatizer
4
5 def prepro_text(target_text):
6     # if b'/b'
7     target_text = re.sub(r"^b[\\'\\"]", '', target_text)
8     target_text = re.sub(r"^[\\w\\s]", '', target_text)
9     target_text = target_text.lower().strip()
10    target_text = target_text.split()
11    target_text = ' '.join([x for x in target_text if x not in stopwords.words('english')])
12    return target_text
```

In [185]:

```

1 # applymap to the whole dataframe
2 #Takes a while to run this because needs to clean every row of data
3
4 merged_data.iloc[:,2] = merged_data.iloc[:,2:].applymap(lambda element:prepr
5 merged_data.head()
6
7 # com_news_df.to_csv('cleaned_Combined_News_DJIA.csv')

```

Out[185]:

	Date	Label	top25_headlines	Top1	Top2	Top3	Top4	T
1988	2008-08-08	0	georgia downs two russian warplanes countries ...	b"Georgia 'downs two Russian warplanes' as cou...	b'BREAKING: Musharraf to be impeached.'	b"Russia Today: Columns of troops roll into So...	b"Russian tanks are moving towards the capital...	b"Afg children re with 'impul U.I
1987	2008-08-11	1	wont america nato help us wont help us help ir...	b"Why wont America and Nato help us? If they w...	b'Bush puts foot down on Georgian conflict'	b"Jewish Georgian minister: Thanks to Israeli ...	b'Georgian army flees in disarray as Russians ...	b'Olyr ope cererr firew 'fak
1986	2008-08-12	0	remember adorable 9yearold sang opening ceremo...	b'Remember that adorable 9-year-old who sang a...	b"Russia 'ends Georgia operation'"	b""If we had no sexual harassment we would hav...	b"Al-Qa'eda is losing support in Iraq because ...	b'Ceasefii Georgia: F Outmaneu ti
1985	2008-08-13	0	us refuses israel weapons attack iran report b...	b' U.S. refuses Israel weapons to attack Iran:...	b"When the president ordered to attack Tskhinv...	b' Israel clears troops who killed Reuters cam...	b'Britain\'s policy of being tough on drugs is...	b"Body c year old fc in trunk; L&
1984	2008-08-14	1	experts admit legalise drugs bwar south osetia...	b'All the experts admit that we should legalis...	b'War in South Osetia - 89 pictures made by a ...	b'Swedish wrestler Ara Abrahamian throws away ...	b"Russia exaggerated the death toll in South O...	b'Missile Killed 9 In Pakistan t

5 rows × 34 columns



Run Model for TFIDF and Embeddings

Import train test split. Use the top 25 headlines combined together as the X variable.

The Y variable is a binary label which shows whether or not the stock market went up or down.

```
In [240]: 1 from sklearn.model_selection import train_test_split
2
3 X = merged_data['top25_headlines']
4 y = merged_data['Label']
5
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

Import the count vectorizer and tensorflow.

```
In [241]: 1 # # Takes a while to load the tensorflow model
2
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.model_selection import train_test_split
5
6 from tqdm.auto import tqdm
7 from sklearn.preprocessing import MaxAbsScaler
8 # Universal Sentence Encoder
9 import tensorflow_hub as hub
10 module_url = "https://tfhub.dev/google/universal-sentence-encoder/4"
11 use = hub.load(module_url)
12 print("module %s loaded" % module_url)
13 import numpy as np
```

module <https://tfhub.dev/google/universal-sentence-encoder/4> (<https://tfhub.dev/google/universal-sentence-encoder/4>) loaded

Import the Tfidf vectorizer.

Create a function to run the model with either TFIDF or embedding.

```

In [242]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2 # return the vectors
3 def text_vector(target_method='tfidf', \
4                 target_list_train=X_train.to_list(),\
5                 target_list_test = X_test.to_list(),\
6                 max_features=None):
7     """
8     type: target_method: string - ("tfidf", "embedding")
9     rtype: list of vectors
10    """
11
12
13    if target_method == "embedding":
14        use_x_train = []
15        use_x_test = []
16        use_x_train = use(target_list_train)
17        use_x_test = use(target_list_test)
18        return use_x_train, use_x_test
19
20    if target_method == "tfidf":
21        vectorizer = TfidfVectorizer(max_features=max_features)
22        tfidf_train = vectorizer.fit_transform(target_list_train)
23        tfidf_test = vectorizer.transform(target_list_test)
24        tfidf_scaler = MaxAbsScaler()
25        tfidf_x_train = tfidf_scaler.fit_transform(tfidf_train)
26        tfidf_x_test = tfidf_scaler.transform(tfidf_test)
27        return tfidf_x_train, tfidf_x_test

```

TFIDF Logistic Regression

```

In [243]: 1 %%time
2
3 tfidf_x_train, tfidf_x_test = text_vector("tfidf", X_train, X_test)
4
5 import pandas as pd
6 import numpy as np
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.metrics import accuracy_score, confusion_matrix
9
10 logmodel = LogisticRegression()
11 logmodel.fit(tfidf_x_train, y_train)
12
13 predictions = logmodel.predict(tfidf_x_test)
14 #print(accuracy_score(y_test, predictions))
15
16 accuracy_tfidf = accuracy_score(y_test, predictions)
17 percentage = "{:.2%}".format(accuracy_tfidf)
18 print("accuracy percentage is", percentage)
19
20

```

accuracy percentage is 51.51%
Wall time: 1.18 s

We got an accuracy below 50% for the TFIDF logistic regression.

Embeddings Logistic Regression

```
In [244]: 1 %%time
2
3 # embedding
4 embedding_x_train, embedding_x_test = text_vector("embedding")
5
6 import pandas as pd
7 import numpy as np
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.metrics import accuracy_score, confusion_matrix
10
11 logmodel = LogisticRegression()
12 logmodel.fit(embedding_x_train, y_train)
13
14 predictions = logmodel.predict(embedding_x_test)
15 #print(accuracy_score(y_test, predictions))
16
17 accuracy_embeddings = accuracy_score(y_test, predictions)
18 percentage = "{:.2%}".format(accuracy_embeddings)
19 print("accuracy percentage is", percentage)
20
21
```

accuracy percentage is 51.01%

Wall time: 3.17 s

We got a slightly better accuracy with embeddings and logistic regression.

MLP classifier

```
In [245]: 1 from sklearn.neural_network import MLPClassifier
2 from sklearn.datasets import make_classification
3
4 clf = MLPClassifier(max_iter = 300).fit(tfidf_x_train, y_train)
5 predictions = clf.predict(tfidf_x_test)
6
7 accuracy_score(y_test, predictions)
8
9 accuracy_stock_mlp = accuracy_score(y_test, predictions)
10 percentage = "{:.2%}".format(accuracy_stock_mlp)
11 print("accuracy percentage is", percentage)
```

accuracy percentage is 47.99%

Text related to stock

Filter the dataset for stock related keywords

Now we want to look at the performance of our models based on only stock related keywords.

Used the column with the top 25 headlines to get the new dataframe with stock related keywords only.

The new dataframe has 553 rows compared to close to 2,000 for the original column.

```
In [226]: 1 stock_related_keywords = "stock|market|feds|bond|bonds|stocks|bull|bear"
2
3 import numpy as np
4
5 temp_df = merged_data.copy()
6
7 merged_data_keywords = temp_df[['top25_headlines', 'Label']]
8
9 merged_data_keywords['Label'] = merged_data_keywords['Label'].astype(str)
10
11 mask = np.column_stack([merged_data_keywords[col].str.contains(stock_related
12 stock_keyword = merged_data_keywords.loc[mask.any(axis=1)]
13
14 stock_keyword['Label'] = stock_keyword['Label'].astype(int)
15
16 #stock_keyword
```

Train test split for stock keyword dataframe.

```
In [235]: 1 from sklearn.model_selection import train_test_split
2
3 X = stock_keyword['top25_headlines']
4 y = stock_keyword.iloc[:,1]
5
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

Run With TFIDF

```
In [236]: 1 tfidf_x_train, tfidf_x_test = text_vector("tfidf", \
2                                             X_train,\
3                                             X_test
4                                             )
```

```
In [233]: 1 from sklearn.linear_model import LogisticRegression
2 logmodel = LogisticRegression()
3 logmodel.fit(tfidf_x_train,y_train)
4 predictions = logmodel.predict(tfidf_x_test)
5
6 accuracy_stock_tfidf = accuracy_score(y_test, predictions)
7 percentage = "{:.2%}".format(accuracy_stock_tfidf)
8 print("accuracy percentage is", percentage)
```

accuracy percentage is 54.55%

Run with MLP

```
In [237]: 1 from sklearn.neural_network import MLPClassifier
2 from sklearn.datasets import make_classification
3
4 clf = MLPClassifier(max_iter = 300).fit(tfidf_x_train, y_train)
5 predictions = clf.predict(tfidf_x_test)
6
7 accuracy_score(y_test, predictions)
8
9 accuracy_stock_mlp = accuracy_score(y_test, predictions)
10 percentage = "{:.2%}".format(accuracy_stock_mlp)
11 print("accuracy percentage is", percentage)
```

accuracy percentage is 56.06%

Sentiment Analysis

Next we will move on to our second research question:

Can the sentiment of reddit daily news headlines be used to predict past stock prices?

Write function to get the sentiment using vader and textblob.

```
In [198]: 1 #uses textblob to get the sentiment scores of the column
2 import nltk
3 nltk.download('vader_lexicon')
4 from textblob import TextBlob
5 import nltk
6 from nltk.sentiment import SentimentIntensityAnalyzer
7 from nltk.sentiment import vader
8 sia = vader.SentimentIntensityAnalyzer()
9
10 def text_sent(target_method, target_list_text):
11     """
12
13     rtype: list of entiment scores
14     """
15     if target_method == 'blob':
16         blob_list = []
17         for title in target_list_text:
18             blob = TextBlob(title)
19             blob_list.append(blob.sentiment.polarity)
20
21         return blob_list
22
23     if target_method == 'NLTK':
24         NLTK_list = []
25
26         for title in target_list_text:
27             sia_polarity = sia.polarity_scores(title)
28             NLTK_list.append(sia_polarity['compound'])
29         return NLTK_list
```

Get the average sentiment for blob and Vader

In [199]:

```

1  # Average sentiment for column by blob
2
3  column_list_blob = []
4
5  for i in range(1,23):
6      #print(i, 'i is')
7      i_str = str(i)
8      Topic = ('Top'+i_str) # 'Top' + 1 = 'Top1'
9      Topic_sent_blob = (Topic + '_sent_' 'blob')
10     topic_list = text_sent("blob",merged_data_sent[Topic])
11     merged_data_sent[Topic_sent_blob] = topic_list
12     column_list_blob.append(Topic_sent_blob)
13
14
15 merged_data_sent_sum_blob = merged_data_sent[column_list_blob].sum(axis=1)
16 merged_data_sent["average_blob_sent"] = merged_data_sent_sum_blob / 22
17
18
19 # Average sentiment by column for NLTK
20
21
22 column_list_nltk = []
23
24 for i in range(1,23):
25     i_str = str(i)
26     Topic = ('Top'+i_str)
27     Topic_sent_nltk = (Topic + '_sent_' 'NLTK')
28     topic_list = text_sent("NLTK",merged_data_sent[Topic])
29     merged_data_sent[Topic_sent_nltk] = topic_list
30     column_list_nltk.append(Topic_sent_nltk)
31
32 merged_data_sent_sum_nltk = merged_data_sent[column_list_nltk].sum(axis=1)
33 merged_data_sent["average_nltk_sent"] = merged_data_sent_sum_nltk / 22
34
35

```

In [200]:

```

1  #merged_data_sent['Top23']

```

Create a DF with just the average blob sent and the average nltk sent. Get the average sent for each day

```
In [201]: 1 label_sent_df= merged_data_sent[['average_nltk_sent', 'average_blob_sent', '
2
3 print('average NLTK sent is' , round(label_sent_df["average_nltk_sent"].mean
4 print('average Blob sent is' , round(label_sent_df["average_blob_sent"].mean
5 label_sent_df
```

average NLTK sent is -0.21

average Blob sent is 0.01

Out[201]:

	average_nltk_sent	average_blob_sent	Label
1988	-0.318659	-0.048722	0
1987	-0.114414	0.030705	1
1986	-0.264577	-0.041955	0
1985	-0.131123	0.005201	0
1984	-0.157518	0.054723	1
...
4	-0.150886	-0.007135	0
3	-0.003755	0.036262	1
2	-0.282536	0.034246	1
1	-0.214800	0.020274	1
0	-0.297809	0.025758	1

1989 rows × 3 columns

Check if positive sentiment for the day is related to a rise in stock prices.

If the Dow Jones went up and the sentiment was positive mark as "correct". If the Dow Jones went up and sentiment was down mark as incorrect.

If the Dow Jones went down and the sentiment was positive mark as "incorrect". If the Dow Jones went down and sentiment was down mark as incorrect.

Do this for both vader and blob.

In [202]:

```
1 nltk_correct = []
2 blob_correct = []
3
4
5 for average_nltk_sent, average_blob_sent, Label in label_sent_df.itertuples():
6     if average_nltk_sent < 0:
7         if Label == 0:
8             correct = 1
9             nltk_correct.append(1)
10        else:
11            correct = 0
12            nltk_correct.append(0)
13    else:
14        #if average_nltk_sent > 0:
15        if Label == 1:
16            correct = 1
17            nltk_correct.append(1)
18        else:
19            correct = 0
20            nltk_correct.append(0)
21
22    if average_blob_sent < 0:
23        if Label == 0:
24            correct = 1
25            blob_correct.append(1)
26        else:
27            correct = 0
28            blob_correct.append(0)
29    else:
30        #if average_nltk_sent > 0:
31        if Label == 1:
32            correct = 1
33            blob_correct.append(1)
34        else:
35            correct = 0
36            blob_correct.append(0)
37
38 label_sent_df['nltk_correct_score'] = nltk_correct
39 label_sent_df['blob_correct_score'] = blob_correct
40
41
42 #label_sent_df
43
44
```

Insert extra columns in DF showing if the blob or vader sent was negative or positive

```
In [203]: 1 import numpy as np
2
3 average_blob_sent_pos = []
4 for x in label_sent_df['average_blob_sent']:
5     if x > 0:
6         average_blob_sent_pos.append(1)
7     else:
8         average_blob_sent_pos.append(0)
9
10 label_sent_df['average_blob_sent_pos'] = average_blob_sent_pos
11
12 average_nltk_sent_pos = []
13 for x in label_sent_df['average_nltk_sent']:
14     if x > 0:
15         average_nltk_sent_pos.append(1)
16     else:
17         average_nltk_sent_pos.append(0)
18
19 label_sent_df['average_nltk_sent_pos'] = average_nltk_sent_pos
20
21
22 label_sent_df['nltk_equals_blob_pos'] = (label_sent_df['average_nltk_sent_po
23
24 #label_sent_df
25
26
```

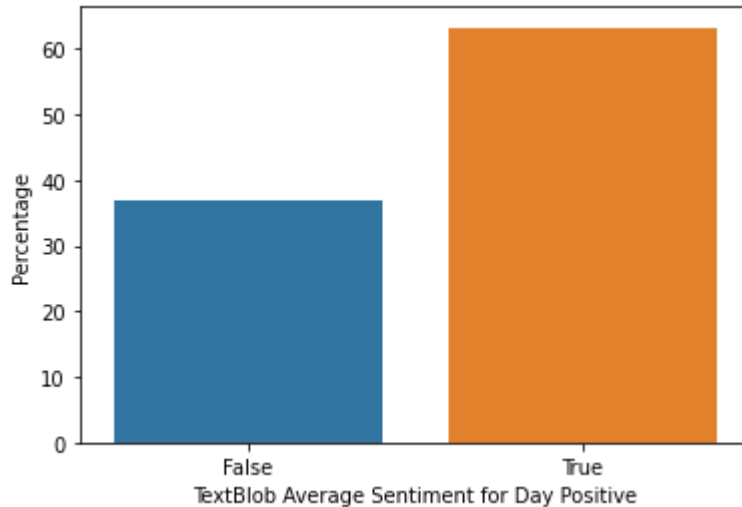
Create plots based of the amount of postives days for both vader and blob sent.

Also compare how many of the vader sentiment matching the blob sentiment.

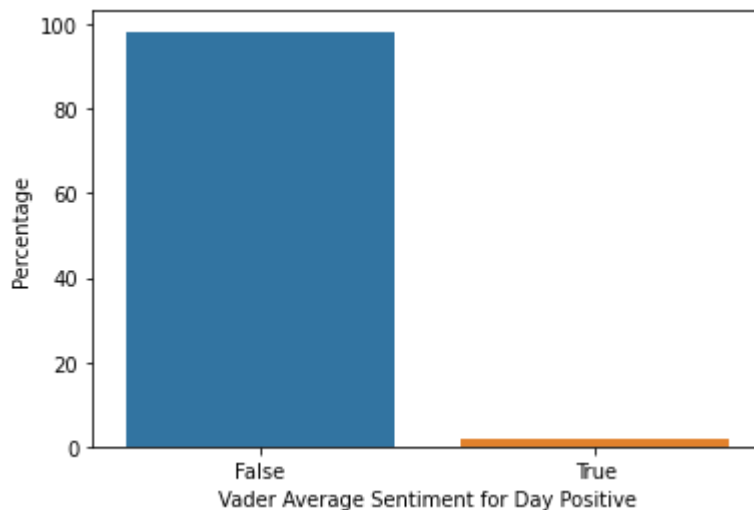
```
In [204]: 1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4
5
6 label_sent_df['nltk_equals_blob_pos'] = label_sent_df.nltk_equals_blob_pos.a
7 label_sent_df['average_blob_sent_pos'] = label_sent_df.average_blob_sent_pos
8 label_sent_df['average_nltk_sent_pos'] = label_sent_df.average_nltk_sent_pos
9 label_sent_df['nltk_correct_score'] = label_sent_df.nltk_correct_score.astype
10 label_sent_df['blob_correct_score'] = label_sent_df.blob_correct_score.astype
```

Vader VS Textblob Average Daily Sentiment

```
In [205]: 1 percentage = lambda i: len(i) / float(len(x)) * 100
2
3 x= label_sent_df['average_blob_sent_pos']
4 ax = sns.barplot(x=x, y=x, estimator=percentage)
5 ax.set(ylabel="Percentage", xlabel = 'TextBlob Average Sentiment for Day Pos
6 plt.show()
7
```



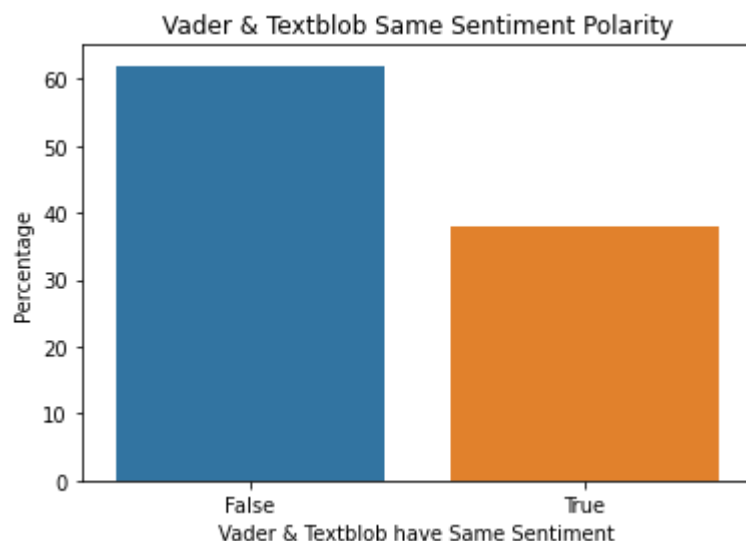
```
In [206]: 1 x= label_sent_df['average_nltk_sent_pos']
2 ax = sns.barplot(x=x, y=x, estimator=percentage)
3 ax.set(ylabel="Percentage", xlabel = 'Vader Average Sentiment for Day Positi
4 plt.show()
```



While Textblob only had 39% of days with an average sentiment score, Vader had over 90% of

days with a negative sentiment score.

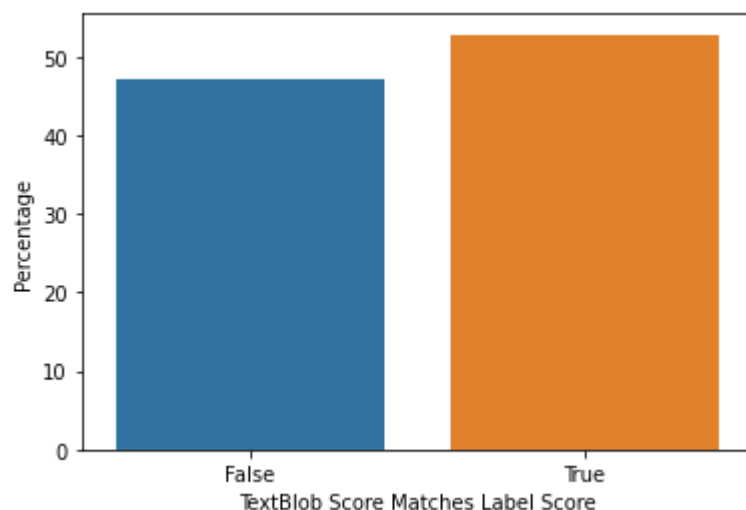
```
In [207]: 1 percentage = lambda i: len(i) / float(len(x)) * 100
          2
          3 x= label_sent_df['nltk_equals_blob_pos']
          4 ax = sns.barplot(x=x, y=x, estimator=percentage)
          5 ax.set(ylabel="Percentage", xlabel = 'Vader & Textblob have Same Sentiment',
          6 plt.show()
```



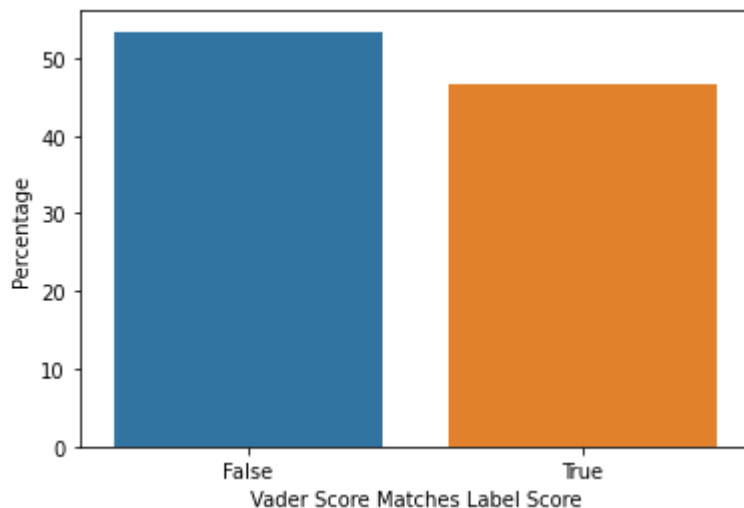
Unsurprisingly, Vader and Textblob only match on polarity (neg and neg or pos and pos) around 40% of the time

Vader and Textblob Match with Dow Jones

```
In [208]: 1 x= label_sent_df['blob_correct_score']
          2 ax = sns.barplot(x=x, y=x, estimator=percentage)
          3 ax.set(ylabel="Percentage", xlabel = 'TextBlob Score Matches Label Score')
          4 plt.show()
          5
```



```
In [209]: 1 x= label_sent_df['nltk_correct_score']  
2 ax = sns.barplot(x=x, y=x, estimator=percentage)  
3 ax.set(ylabel="Percentage", xlabel = 'Vader Score Matches Label Score')  
4 plt.show()
```



Textblob does a slightly better job of matching the Dow Jones average, however both Textblob and Vader match right around 50% of the time.

Sentiment With Decision Tree and Logistic Regression Models

Logistic Regression Model

Perform a logistic regression model using sentiment to predict if stocks went up or down.

Use text blob and NLTK as the X and the Dow Jones label of whether Dow Jones went down or up as the y.

Textblob

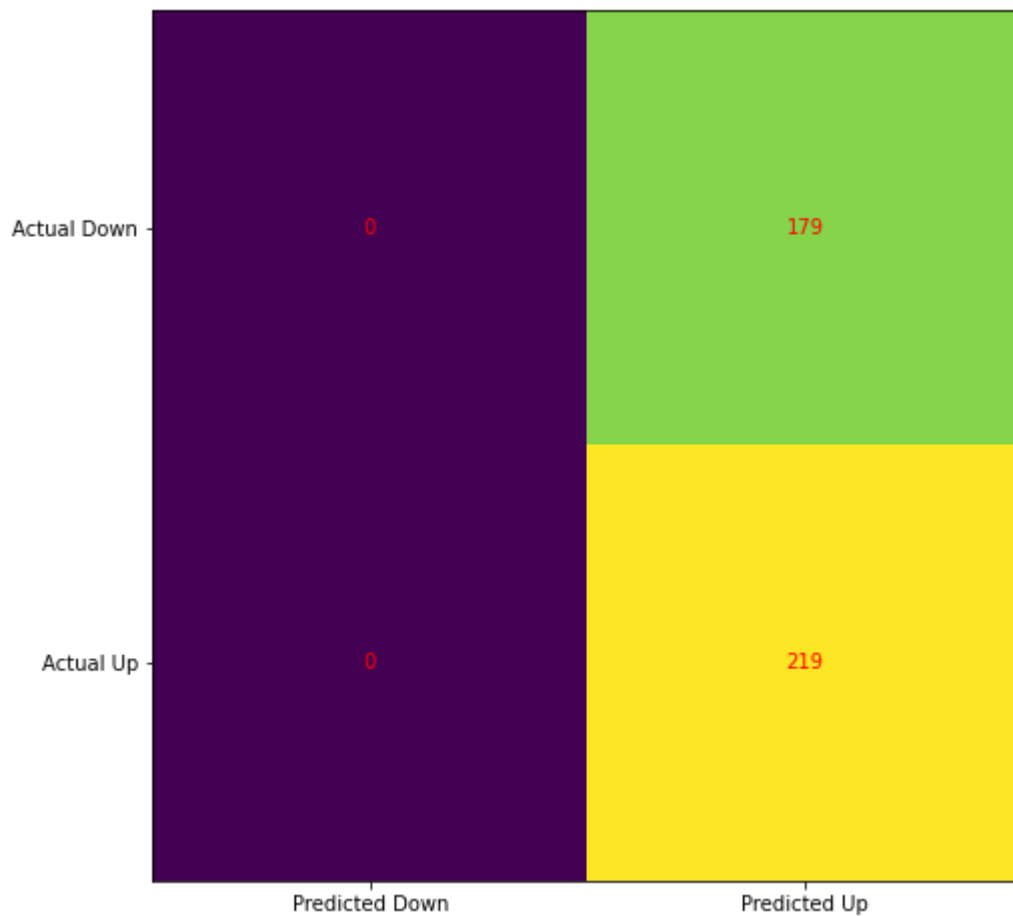
```
In [210]: 1 from sklearn.model_selection import train_test_split
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.metrics import classification_report, confusion_matrix
6
7 model = LogisticRegression(solver='liblinear')
8
9
10 X = merged_data_sent[['average_blob_sent']]
11 y = merged_data_sent['Label']
12
13
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```

In [211]: 1 from sklearn.linear_model import LogisticRegression
2 logmodel = LogisticRegression()
3 logmodel.fit(X_train,y_train)
4 predictions_blob = logmodel.predict(X_test)
5
6 from sklearn.metrics import classification_report
7 print(classification_report(y_test,predictions_blob))
8
9 cm = confusion_matrix(y_test, predictions_blob)
10
11 fig, ax = plt.subplots(figsize=(8, 8))
12 ax.imshow(cm)
13 ax.grid(False)
14 ax.xaxis.set(ticks=(0, 1), ticklabels=('Predicted Down', 'Predicted Up'))
15 ax.yaxis.set(ticks=(0, 1), ticklabels=('Actual Down', 'Actual Up'))
16 ax.set_ylim(1.5, -0.5)
17 for i in range(2):
18     for j in range(2):
19         ax.text(j, i, cm[i, j], ha='center', va='center', color='red')
20 plt.show()

```

	precision	recall	f1-score	support
0	0.00	0.00	0.00	179
1	0.55	1.00	0.71	219
accuracy			0.55	398
macro avg	0.28	0.50	0.35	398
weighted avg	0.30	0.55	0.39	398



Vader

```
In [212]: 1 #need to split into blob and NLTK and not combine together.
2
3
4
5 from sklearn.model_selection import train_test_split
6 import matplotlib.pyplot as plt
7 import numpy as np
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.metrics import classification_report, confusion_matrix
10
11 model = LogisticRegression(solver='liblinear', random_state=0)
12
13
14 X = merged_data_sent[['average_nltk_sent']]
15 y = merged_data_sent['Label']
16
17
18 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

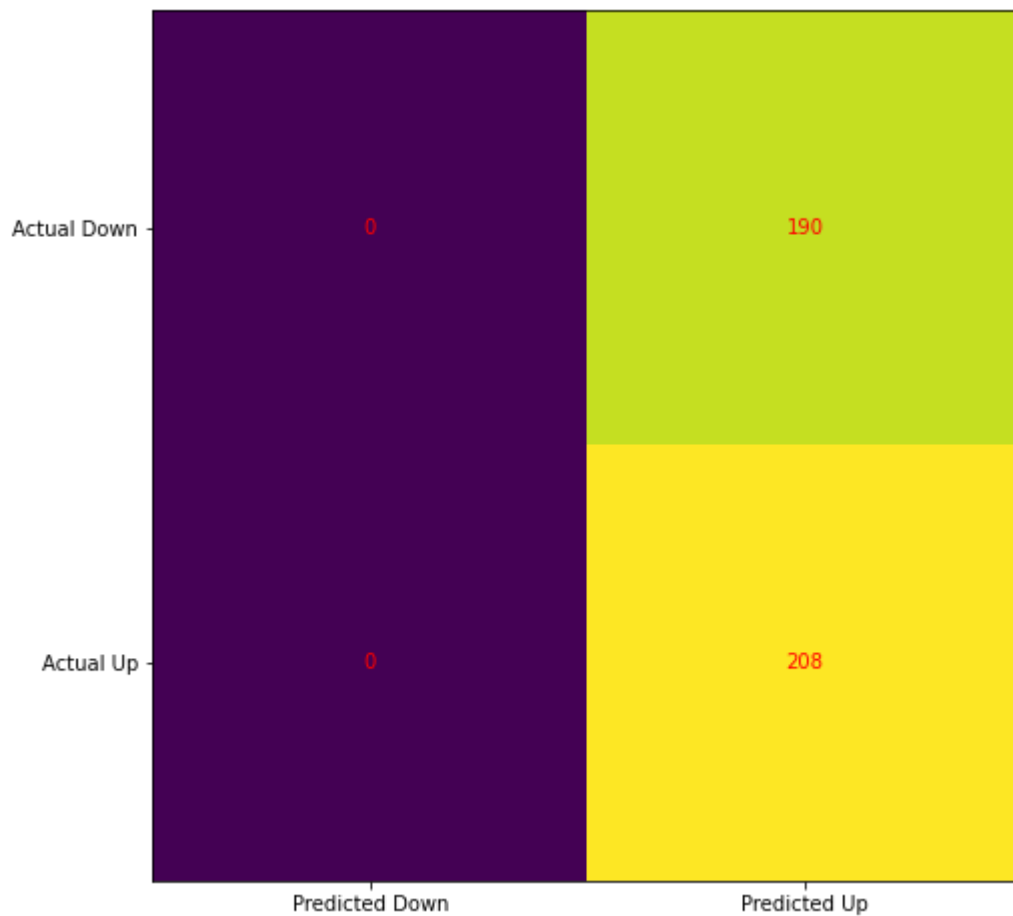
In [213]:

```

1 from sklearn.linear_model import LogisticRegression
2 logmodel = LogisticRegression()
3 logmodel.fit(X_train,y_train)
4
5 predictions_nltk = logmodel.predict(X_test)
6
7 from sklearn.metrics import classification_report
8 print(classification_report(y_test,predictions_nltk))
9
10 cm = confusion_matrix(y_test, predictions_nltk)
11
12 fig, ax = plt.subplots(figsize=(8, 8))
13 ax.imshow(cm)
14 ax.grid(False)
15 ax.xaxis.set(ticks=(0, 1), ticklabels=('Predicted Down', 'Predicted Up'))
16 ax.yaxis.set(ticks=(0, 1), ticklabels=('Actual Down', 'Actual Up'))
17 ax.set_ylim(1.5, -0.5)
18 for i in range(2):
19     for j in range(2):
20         ax.text(j, i, cm[i, j], ha='center', va='center', color='red')
21 plt.show()

```

	precision	recall	f1-score	support
0	0.00	0.00	0.00	190
1	0.52	1.00	0.69	208
accuracy			0.52	398
macro avg	0.26	0.50	0.34	398
weighted avg	0.27	0.52	0.36	398



The vader logistic regression model was slightly more accurate than the textblob.

Decision Tree Classifier

TextBlob

```
In [214]: 1 from sklearn.tree import DecisionTreeClassifier
2
3 classifier_tree = DecisionTreeClassifier()
4
5 X_blob = merged_data_sent[['average_blob_sent']]
6 y_blob = merged_data_sent['Label']
7
8
9 X_train, X_test, y_train, y_test = train_test_split(X_blob, y_blob, test_siz
```

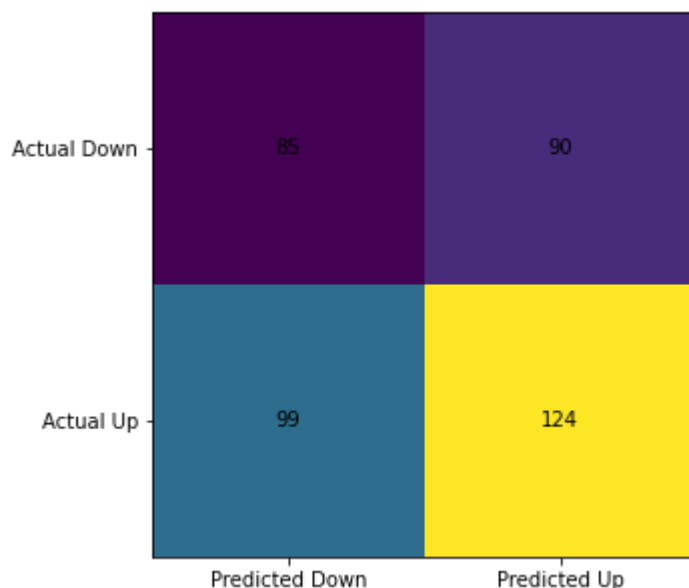
```

In [215]: 1 y_predict_blob = classifier_tree.fit(X_train, y_train).predict(X_test)
          2
          3 print("Decision Tree Classifier using text blob \n",classification_report(y_
          4
          5
          6 cm = confusion_matrix(y_test, y_predict_blob)
          7
          8 fig, ax = plt.subplots(figsize=(5, 5))
          9 ax.imshow(cm)
         10 ax.grid(False)
         11 ax.xaxis.set(ticks=(0, 1), ticklabels=('Predicted Down', 'Predicted Up'))
         12 ax.yaxis.set(ticks=(0, 1), ticklabels=('Actual Down', 'Actual Up'))
         13 ax.set_ylim(1.5, -0.5)
         14 for i in range(2):
         15     for j in range(2):
         16         ax.text(j, i, cm[i, j], ha='center', va='center', color='black')
         17 plt.show()

```

Decision Tree Classifier using text blob

	precision	recall	f1-score	support
0	0.46	0.49	0.47	175
1	0.58	0.56	0.57	223
accuracy			0.53	398
macro avg	0.52	0.52	0.52	398
weighted avg	0.53	0.53	0.53	398



Vader

```
In [216]: 1 from sklearn.tree import DecisionTreeClassifier
          2
          3 classifier_tree = DecisionTreeClassifier()
          4
          5 X_nltk = merged_data_sent[['average_nltk_sent']]
          6 y_nltk = merged_data_sent['Label']
          7
          8
          9 X_train, X_test, y_train, y_test = train_test_split(X_nltk, y_nltk, test_siz
```

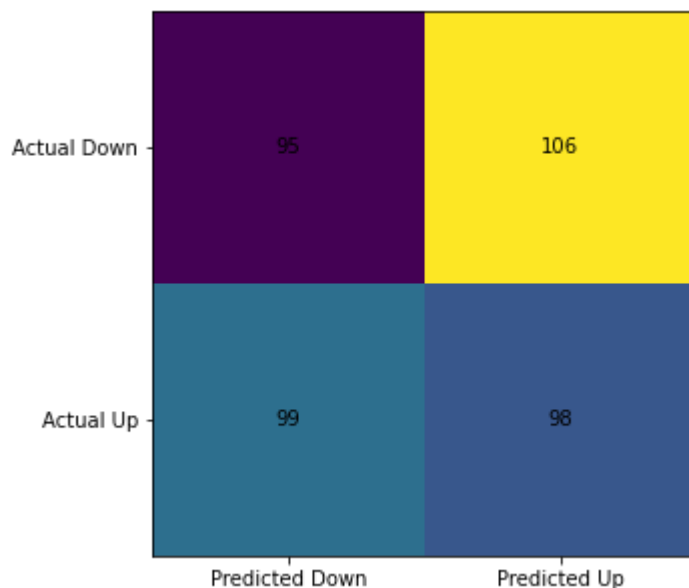
```

In [217]: 1 y_predict_nltk = classifier_tree.fit(X_train, y_train).predict(X_test)
          2
          3 print("Decision Tree Classifier using NLTK \n", classification_report(y_test
          4
          5
          6 cm = confusion_matrix(y_test, y_predict_nltk)
          7
          8 fig, ax = plt.subplots(figsize=(5, 5))
          9 ax.imshow(cm)
         10 ax.grid(False)
         11 ax.xaxis.set(ticks=(0, 1), ticklabels=('Predicted Down', 'Predicted Up'))
         12 ax.yaxis.set(ticks=(0, 1), ticklabels=('Actual Down', 'Actual Up'))
         13 ax.set_ylim(1.5, -0.5)
         14 for i in range(2):
         15     for j in range(2):
         16         ax.text(j, i, cm[i, j], ha='center', va='center', color='black')
         17 plt.show()

```

Decision Tree Classifier using NLTK

	precision	recall	f1-score	support
0	0.49	0.47	0.48	201
1	0.48	0.50	0.49	197
accuracy			0.48	398
macro avg	0.49	0.49	0.48	398
weighted avg	0.49	0.48	0.48	398



For the decision tree, both vader and textblob had similar accuracy results.

In []:

1

In []:

1

Overall Conclusions

TFIDF and Embeddings

- In general, our models do not predict the trend of the Dow Jones Industrial Average of each day very well.
- Even though we achieved 51% and 57% accuracy scores, they are still too low to say with any certainty that they are good at predicting the market. With related keywords such as “stock” or “market”, we have increased accuracy by over 5% for our logistic regression model.
- We believe that our input variables are not noisy, because we were able to increase our accuracy for machine learning models by adding filters and selecting more related variables out of the original data sets, such as selecting the texts with the key words, etc.

Sentiment

- Textblob and Vader produced extremely different average sentiment scores.
- Both Textblob and Vader have average accuracies around 53%, with the best average accuracy at 54% and logistic regression with slightly higher accuracies.
- Our decision tree model has a realistic distribution than the logistic regression model as it is a lot less skewed.