

Sentiment Analysis of Twitter Data for predicting Stock Price

Introduction: This project involves evaluation TESLA's social media data and uses sentiment classification to estimate its future stock trend. We employed sentiment analysis and machine learning principles to estimate the impact of "public sentiment" on "market fluctuations"

Import Required Python Packages

```
import pandas as pd
import numpy as np
import requests
from flask import Response
import re

from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.ensemble import RandomForestClassifier

from matplotlib import rcParams
rcParams['figure.figsize'] = 12, 6
import matplotlib.pyplot as plt
plt.style.use('ggplot')

from sklearn.metrics import *
import warnings
warnings.filterwarnings("ignore")

from Sweet.sweet import scrape
from Sweet.user import get_user_information, get_users_following, get_users_followers
```

```
sentiment_model = flair.models.TextClassifier.load('en-sentiment')

2021-11-10 02:41:134,974 loading file C:\Users\fhahad\flair\models\sentiment-en-mix-distillbert_4.pt
```

Scweet Python Scraper

```
# Python API allows to scrape only 7 days of data in free tier, So we managed to use Scweet tweets scraper E
#data = scrape(words=['tesla'], since="2020-10-31", until="2021-10-31",interval=1,
# headless=True, display_type="Top")

# data saved to csv file after scraping
#data.to_csv('scraped_data.csv', index = False)

data = pd.read_csv('scraped_data.csv')
data.head()
```

	UserScreenName	UserName	Timestamp	Text	Embedded text	Emojis	Comments	Likes	Retw
0	East Bay Tesla Club Fremont	@TeslaOwnersBay	2020-10-31T23:10:55.000Z	East Bay Tesla Club Fremont	Happy Halloween all you @Tesla!nFreaks!n...	👻👻👻		5	40
1	Kim Paquette	@kimpaquette	2021-10-24:05.000Z	Kim Paquette@kimpaquette@n Nov 2020	Guess what I am for Halloween!n#FSDBeta #Tes...	👻👻👻		150	46
2	Dogs of Tesla	@dogsofTesla	2020-10-31T21:58:30.000Z	Dogs of Tesla@dogsofTesla@n Nov 2020	FRANKIE nHappy Halloween!n@Tesla...	👻👻👻		10	21
3	Luke	@tweetofbfc	2020-10-31T22:54:04.000Z	Luke@tweetofbfc@n Nov 2020	Major vibes anyone know whose this is? #tesla...	👻		2	4
4	Randy 'DUO' Mongenel	@RandyMongenel	2020-10-31T20:22:41.000Z	Mongenel@RandyMongenel@n Nov...	One of the more interesting decals I've seen o...	NaN		2	NaN

Select only the Date and Text of speech from scraped data

```
import datetime
data['date'] = pd.to_datetime(data['Timestamp']).dt.date
df = data[['date', 'Embedded text']]
df.head(20)
```

	date	Embedded text
0	2020-10-31	Happy Halloween all you @Tesla!nFreaks!n...
1	2020-10-31	Guess what I am for Halloween!n#FSDBeta #Tes...
2	2020-10-31	FRANKIE nHappy Halloween!n@Tesla...
3	2020-10-31	Major vibes anyone know whose this is? #tesla...
4	2020-10-31	One of the more interesting decals I've seen o...
5	2020-10-31	Just realized that even with the dozen or so F...
6	2020-10-31	How did you get a hold of that Chester? Thank...
7	2020-10-31	Our #socialdistanced #rickortreat solution...
8	2020-10-31	zombie apocalypse happen i'm stealing a tesla...
9	2020-10-31	Demand problem Natick Massachusetts. You know...
10	2020-10-31	What kind of science do you think is the futur...
11	2020-10-31	Hiring Bo Pelini to run an SEC defense in 202...
12	2020-10-31	"I do not think you can name many great inven...
13	2020-10-31	Patently waiting for our new and final @Tesla...
14	2020-10-31	The all-new, fully electric #Polstar2 is the ...
15	2020-10-31	What is that @Tesla@n @btp.pl@n @HSBC@n @...
16	2020-10-31	My kid just pointed out to me that when the @n...
17	2020-10-31	Can't wait for @Tesla@n #FSDbeta to be rolle...
18	2020-10-31	Replying to @n@dogsofTesla@n @n@elonmusk@n an...
19	2020-10-31	Replying to @n@dogsofTesla@n @n@elonmusk@n an...

Combine the tweets for each date in same rows.

```
df['Embedded text'] = df.groupby('date')[['Embedded text']].transform(lambda x : ' '.join(x))
x = df.drop_duplicates().reset_index().drop('index', axis = 1)
df.head()
```

	date	Embedded text
0	2020-10-31	Happy Halloween all you @Tesla!nFreaks!n...
1	2020-11-01	Is this the first Tesla in the country?n1@n...
2	2020-11-02	Tesla Model Y is here! And quality is pretty d...
3	2020-11-03	I like that the Tesla maps now identify inters...
4	2020-11-04	Hey Tesla take me to Kinoo@nAutopilot.n@n...

Sentiment Analysis

- Flair
- Analyzing Tesla Tweets

```
# function to clean tweet text by removing links, special characters using simple regex statements
def clean_tweet(tweet):
    """Utility to clean a tweet"""
    return ' '.join(re.sub('([@&#%*~&#x20-9+])|((\r|0-9A-Z-a-z \t)|(\n+/\r+/$))', "", tweet).split())
```

Notes: We'll be using the flair library's pre-trained sentiment analysis model. This model divides the text into character-level tokens and makes predictions using the BiLSTM model. Working at the character level (rather than the word level) has the advantage of allowing the network to attribute a feeling to words it has never encountered before.

```
proba = []
sentiments = []

for tweet in df['Embedded text'].to_list():
    # make prediction
    tweet = clean_tweet(tweet)
    sentence = flair.data.Sentence(tweet)
    sentiment_model.predict(sentence)
    #extract sentiment prediction
    probs.append(sentence.labels[0].score) # numerical score 0-1
    sentiments.append(sentence.labels[0].value) # "POSITIVE" or "NEGATIVE"
```

```
# add probability and sentiment predictions to tweets dataframe
df['probability'] = probs
df['sentiment'] = sentiments
```

```
tweet_prob = df[['date', 'probability', 'sentiment']]
tweet_prob['date'] = pd.to_datetime(tweet_prob['date']) # convert string date object to datetime object
```

Use Textblob to calculate polarity parameter from the tweet text

```
from textblob import TextBlob
from nltk.sentiment.vader import SentimentIntensityAnalyzer

polarity = 0
tweet_polarity = []

# print(tweet.text)
for tweet in df['Embedded text'].to_list():
    analysis = TextBlob(tweet)
    score = SentimentIntensityAnalyzer().polarity_score(tweet)
    polarity = analysis.sentiment.polarity
    tweet_polarity.append(polarity)
```

Inference: We can observe from the above results flair model has created the positive and negative sentiments for compiled tweets of each day of last one year.

```
# The yahoo finance API is not working fine to scrape the data, so we directly
# downloaded one year of data from https://finance.yahoo.com/quote/TSLA/
tesla = pd.read_csv('tesla.csv')
tesla['date'] = pd.to_datetime(tesla['Date']) # convert string date object to datetime object
tesla
```

	Date	Open	High	Low	Close	Adj Close	Volume
0	2020-02-11	394.000000	406.980011	392.299988	400.510010	400.510010	29021100
1	2020-03-11	409.730011	427.769989	406.690002	423.899994	423.899994	34351700
2	2020-04-11	430.619995	435.399994	417.100006	420.980011	420.980011	32143100
3	2020-05-11	428.299988	440.000000	424.000000	438.089996	438.089996	28414500
4	2020-06-11	436.100006	436.570007	424.279999	429.950012	429.950012	21706000
...
246	2021-10-25	950.530029	1045.020020	944.200012	1024.859985	1024.859985	62852100
247	2021-10-26	1024.689941	1094.939941	1001.440002	1018.429993	1018.429993	62415000
248	2021-10-27	1039.660034	1070.880005	1030.780029	1037.859985	1037.859985	38526500
249	2021-10-28	1068.310059	1081.000000	1054.199951	1077.040039	1077.040039	27213200
250	2021-10-29	1081.859985	1115.209961	1073.209961	1114.000000	1114.000000	29918400

251 rows x 7 columns

```
df_merged = pd.merge(left=tweet_prob, right=tesla, left_on='date', right_on='Date')
df_merged
```

	date	probability	sentiment	tweet_polarity	Date	Open	High	Low	Close	Adj Close	Volume
0	2020-11-11	0.821078	POSITIVE	0.150447	2020-11-11	416.450012	418.700012	410.579987	417.130005	417.130005	17357700
1	2020-11-12	0.942596	NEGATIVE	0.217726	2020-11-12	615.010010	624.000000	596.799988	609.989990	609.989990	46475000
2	2020-11-13	0.999725	NEGATIVE	0.092293	2020-11-13	410.850006	412.529999	401.660004	408.500000	408.500000	19771100
3	2020-11-16	0.998980	POSITIVE	0.200596	2020-11-16	408.929993	412.450012	404.089996	408.089996	408.089996	26838600
4	2020-11-17	0.998251	POSITIVE	0.192465	2020-11-17	460.170013	462.000000	433.010010	441.609985	441.609985	61188300
...
211	2021-10-25	0.998716	POSITIVE	0.190844	2021-10-25	950.530029	1045.020020	944.200012	1024.859985	1024.859985	62852100
212	2021-10-26	0.999432	NEGATIVE	0.124890	2021-10-26	1024.689941	1094.939941	1001.440002	1018.429993	1018.429993	62415000
213	2021-10-27	0.999017	NEGATIVE	0.085150	2021-10-27	1039.660034	1070.880005	1030.780029	1037.859985	1037.859985	38526500
214	2021-10-28	0.534431	POSITIVE	0.192312	2021-10-28	1068.310059	1081.000000	1054.199951	1077.040039	1077.040039	27213200
215	2021-10-29	0.885732	POSITIVE	0.241323	2021-10-29	1081.859985	1115.209961	1073.209961	1114.000000	1114.000000	29918400

216 rows x 11 columns

```
#stock Trend is calculated by taking the difference between today's price and the previous day's
#price and converting it to a '0' for a decline and a '1' for an increase in stock price.
df_merged['Price Diff'] = df_merged['Adj Close'].diff()
df_merged.dropna(inplace = True)
df_merged['Trend'] = np.where(df_merged['Price Diff'] > 0 , 1, 0)
df_merged.head()
```

	date	probability	sentiment	tweet_polarity	Date	Open	High	Low	Close	Adj Close	Volume	Price Diff
1	2020-11-12	0.942596	NEGATIVE	0.217726	2020-11-12	615.010010	624.000000	596.799988	609.989990	609.989990	46475000	192.859985
2	2020-11-13	0.999725	NEGATIVE	0.092293	2020-11-13	410.850006	412.529999	401.660004	408.500000	408.500000	19771100	-201.489990
3	2020-11-16	0.998980	POSITIVE	0.200596	2020-11-16	408.929993	412.450012	404.089996	408.089996	408.089996	26838600	-0.410004
4	2020-11-17	0.998251	POSITIVE	0.192465	2020-11-17	460.170013	462.000000	433.010010	441.609985	441.609985	61188300	33.519989
5	2020-11-18	0.997662	NEGATIVE	0.201802	2020-11-18	448.350006	496.000000	443.500000	486.640015	486.640015	78044000	45.030030

Create separate binary variables for Sentiment features

```
# encode each categorical variable in the loop
df_e = pd.get_dummies(df_merged['sentiment'], prefix = 'sentiment', drop_first = False)
# combine the insurance dataset with the encoded categorical variable
df_merged = pd.concat([df_merged, df_e], axis = 1)
```

```
# drop the original non-encoded categorical variable from the dataset
del df_merged['sentiment']
df_merged = df_merged.set_index('date')
```

```
df_merged.head()
```

	probability	tweet_polarity	Date	Open	High	Low	Close	Adj Close	Volume	Price Diff	Trend	sentim
date												
2020-11-12	0.942596	0.217726	2020-11-12	615.010010	624.000000	596.799988	609.989990	609.989990	46475000	192.859985		1
2020-11-13	0.999725	0.092293	2020-11-13	410.850006	412.529999	401.660004	408.500000	408.500000	19771100	-201.489990		0
2020-11-16	0.998980	0.200596	2020-11-16	408.929993	412.450012	404.089996	408.089996	408.089996	26838600	-0.410004		0
2020-11-17	0.998251	0.192465	2020-11-17	460.170013	462.000000	433.010010	441.609985	441.609985	61188300	33.519989		1
2020-11-18	0.997662	0.201802	2020-11-18	448.350006	496.000000	443.500000	486.640015	486.640015	78044000	45.030030		1

Predicting TESLA Stock TREND with Random Forest Classifier Model

```
# Select independent and dependent variables
X = df_merged[['Adj Close', 'Volume', 'tweet_polarity', 'sentiment_NEGATIVE', 'sentiment_POSITIVE']]
y = df_merged['Trend']
```

```
# Scale the X data onto same scale
from sklearn import preprocessing
X = preprocessing.StandardScaler().fit(X).transform(X.astype(float))
X[5]
```

```
array([[0.87246471, 0.85095057, 0.97004924, 0.61498384, -0.61498384],
       [-2.7494569 , -0.54696907, -1.67908558, 0.61498384, -0.61498384],
       [-2.75327632, -0.17700516, 0.60826108, -1.62605898, 1.62605898],
       [-2.44101983, 1.62110589, 0.43653402, -1.62605898, 1.62605898],
       [-2.02153886, 2.50345474, 0.63374177, -0.61498384, -0.61498384]])
```

Split the data into Train and Test sets

```
X_train, X_test, y_train, y_test = train_test_split(X,y,test_size = 0.25, random_state = 40)
```

Create a Random Forest Classifier Model

```
# Create RFClassifier model
rf_model = RandomForestClassifier(n_estimators=500, random_state=78)
# Fit the model
rf_model = rf_model.fit(X_train, y_train)
```

```
# Make predictions
predictions = rf_model.predict(X_test)
pd.DataFrame({"Prediction": predictions, "Actual": y_test}).head(20)
```

```
# Generate accuracy score for predictions using y_test
acc_score = accuracy_score(y_test, predictions)
print(f"Accuracy Score : {acc_score}")
```

```
Accuracy Score : 0.5185185185185185
```

```
cm = confusion_matrix(y_test, predictions)
cm_df = pd.DataFrame(
    cm, index=["Actual 0", "Actual 1"],
    columns=["Predicted 0", "Predicted 1"]
)
display(cm_df)
```

	Predicted 0	Predicted 1
Actual 0	8	21
Actual 1	5	20

```
# Generating classification report
print("Classification Report")
print(classification_report(y_test, predictions))
```

stocks_1.head()				
	Real	Predicted		
date				
2021-08-10	653.380005	650.765625		
2021-08-13	655.289978	657.669189		
2021-08-16	846.640015	837.144958		
2021-08-17	708.489990	747.709839		
2021-08-18	677.349976	671.159851		

Create a Gradient Booster Model

```
# Choosing learning rate
model = GradientBoostingClassifier(n_estimators=20,
                                   learning_rate=0.05,
                                   max_features=2,
                                   max_depth=3,
                                   random_state=0)
```

```
model.fit(X_train, y_train)
GB_pred = model.predict(X_test)
```

```
# Scoring the model
print("Accuracy score (training): {0:3f}".format(
    model.score(
        X_train,
        y_train)))
```

```
print("Accuracy score (validation): {0:3f}".format(
    model.score(
        X_test,
        y_test)))
```

```
print()
```

```
Accuracy score (training): 0.801
Accuracy score (validation): 0.426
```

Evaluation of a Model

Using sklearn to build the confusion matrix and generate the classification report, evaluating the model's findings.

```
cm = confusion_matrix(y_test, GB_pred)
cm_df = pd.DataFrame(
    cm, index=["Actual 0", "Actual 1"],
    columns=["Predicted 0", "Predicted 1"]
)
display(cm_df)
```

	Predicted 0	Predicted 1
Actual 0	1	28
Actual 1	3	22

```
# Generating classification report
print("Classification Report")
print(classification_report(y_test, GB_pred))
```

```
<AxesSubplot:title=('center','Real vs Predicted values of TESLA'), xlabel='date'>
```

The chart displays two time series: 'Real' (red line) and 'Predicted' (blue line). The y-axis represents values ranging from 750 to 850. The x-axis is labeled 'date'. The 'Real' series shows several sharp peaks, notably around 840, 860, and 880. The 'Predicted' series follows the general trend of the 'Real' series but tends to miss the highest peaks, reaching a maximum of approximately 800.

Random Forest Regressor Model - TESLA Stock Price Prediction

In this approach the tweet polarity, price differences the adjusted closing prices were used to predict the future movement in Apple stock price. The adjusted closing prices data was structured such that previous time steps were used as input variables and the next time step as the output variable.

Select required data for regression analysis variables

```
df2 = df_merged[['Adj Close', 'Price Diff', 'sentiment_POSITIVE', 'Volume', 'tweet_polarity']]
df2.head()
```

	Adj Close	Price Diff	sentiment_POSITIVE	Volume	tweet_polarity
date					
2020-11-12	609.989990	192.859985	0	46475000	0.217726
2020-11-13	408.500000	-201.489990	0	19771100	0.092293
2020-11-16	408.089996	-0.410004	1	26838600	0.200596
2020-11-17	441.609985	33.51998			

