# Determining the Sentiment of Financial News

The news database here will train the Naive Bayes, then RandomForrest For deploying I'd recommend using the NewsAPI code shared and tag the sentiment via the trained NB.

## Constructing a Naive Bayes Classifier

- Load dataset
- Vectorize data
- Split data (80/20, train test, random_state=0 so as to allow reproducability)
- Initialize the NB classifer and fit
- Predict and measure accuracy

In [9]:
```python
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_c
from sklearn.feature_extraction.text import CountVectorizer

news_pd = pd.read_csv("./news_with_sentiment.csv")
news_pd = news_pd[:2000] # 28,000 rows will use more RAM than is av

cv = CountVectorizer() # Convert text data to a vector as that is r
X = cv.fit_transform(news_pd['text']).toarray()
y = news_pd['sentiment'] # y = the variable we are trying to predic
```

In [10]:
```python
# Split train and test data (80/20)
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size

# Initialize the Gaussian Naive Bayes Classifier, then fit the data
from sklearn.naive_bayes import GaussianNB
classifier = GaussianNB()
classifier.fit(X_train, y_train)
```

Out[10]: GaussianNB(priors=None)

In [11]:
```python
# Predict sentiment of our test data
y_pred = classifier.predict(X_test)

from sklearn.metrics import accuracy_score
score = accuracy_score(y_test, y_pred)
```

And now we can view the accuracy:

In [12]:    1   print(score)

0.685

| 1 | Roughly 68% accuracy. Not exactly stellar, if you reduce the dataset further you end up with higher accuracy which is interesting. |

In [14]:
```
 1   news_pd = pd.read_csv("./news_with_sentiment.csv")
 2   news_pd = news_pd[:1000] # 28,000 rows will use more RAM than is av
 3
 4   cv = CountVectorizer()
 5   X = cv.fit_transform(news_pd['text']).toarray()
 6   y = news_pd['sentiment']
 7
 8   X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
 9
10   classifier = GaussianNB()
11   classifier.fit(X_train, y_train)
12
13   y_pred = classifier.predict(X_test)
14
15   score = accuracy_score(y_test, y_pred)
16
17   print(score)
```

0.775

77.5% accuracy on a 1000 row dataset with an 80/20 split.

After research, Naive Bayes appears to be better with smaller datasets but perhaps we can improve:

## To improve on our Naive Bayes we can now try a Random Forest:

- Load dataset
- Remove stopwords, min_df=7 means the data is irrelevant if used in more than 7 documents, max_df of 0.8 means it also is irrelevant if used in more than 80% of documents
- Vectorize data (max_features is the max number of WORDS in Vector form that will influence the sentiment)
- Split data (80/20, train test, random_state=0 so as to allow reproducability)
- Initialize the Random Forest classifer and fit
- Predict and measure accuracy

```
In [15]:  1  # Read in 20,000 headlines
          2  news_pd = pd.read_csv("./news_with_sentiment.csv")
          3  news_pd = news_pd[:20000] # 28,000 rows will use more RAM than is a
          4  y = news_pd['sentiment']
```

```
In [16]:  1  from nltk.corpus import stopwords
          2  from sklearn.feature_extraction.text import TfidfVectorizer
          3
          4  # Remove stopwords and vectorize the dataset
          5  #TfidVectorizer converts a collection of raw documents to a matrix
          6  vectorizer = TfidfVectorizer(max_features=2500, min_df=7, max_df=0.
          7  processed_features = vectorizer.fit_transform(news_pd['text']).toar
```

```
In [17]:  1  # 80/20 data split
          2  from sklearn.model_selection import train_test_split
          3  X_train, X_test, y_train, y_test = train_test_split(processed_featu
          4
          5  # Fit our model with split data, starting with 450 estimators (450
          6  from sklearn.ensemble import RandomForestClassifier
          7
          8  text_classifier = RandomForestClassifier(n_estimators=450, random_s
          9  text_classifier.fit(X_train, y_train)
```

```
/home/nbuser/anaconda3_420/lib/python3.5/site-packages/sklearn/ensembl
e/weight_boosting.py:29: DeprecationWarning: numpy.core.umath_tests is
an internal NumPy module and should not be imported. It will be remove
d in a future NumPy release.
  from numpy.core.umath_tests import inner1d
```

```
Out[17]: RandomForestClassifier(bootstrap=True, class_weight=None, criterion='g
ini',
            max_depth=None, max_features='auto', max_leaf_nodes=None,
            min_impurity_decrease=0.0, min_impurity_split=None,
            min_samples_leaf=1, min_samples_split=2,
            min_weight_fraction_leaf=0.0, n_estimators=450, n_jobs=1,
            oob_score=False, random_state=0, verbose=0, warm_start=Fal
se)
```

```
In [18]:  1  # Predicting the sentiment of our test data
          2  predictions = text_classifier.predict(X_test)
          3
          4
          5  # Checking our accuracy
          6  from sklearn.metrics import accuracy_score
          7  print(accuracy_score(y_test, predictions))
```

```
0.93575
```

93.57% accuracy

# Hyperparameter Tuning:

- Choose a set of trees we want to test
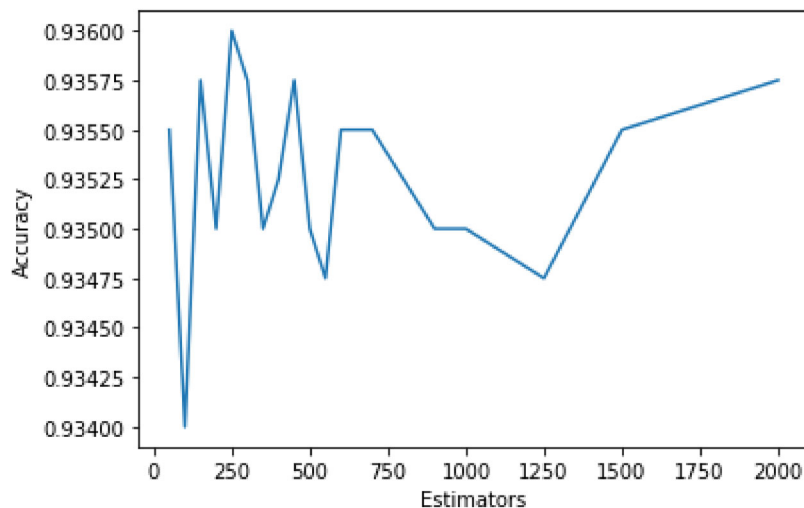- Train the model with n trees, store accuracy

```
In [ ]:     1  from sklearn.ensemble import RandomForestRegressor
            2
            3  estimators = [50, 100, 150, 200, 250, 300, 350, 400, 450, 500, 550,
            4  accuracy = []
            5
            6  for estimator_num in estimators:
            7      # Fit and predict
            8      text_classifier = RandomForestClassifier(n_estimators=estimator
            9      text_classifier.fit(X_train, y_train)
           10      predictions = text_classifier.predict(X_test)
           11
           12      # Store accuracy
           13      from sklearn.metrics import accuracy_score
           14      accuracy.append(accuracy_score(y_test, predictions))
           15
           16
           17  # Graph reported accuracy of various sets of estimators
           18  import matplotlib.pyplot as plt
           19
           20  plt.plot(estimators, accuracy)
           21  plt.ylabel('Accuracy')
           22  plt.xlabel('Estimators')
           23  plt.show()
           24
           25  print(estimators)
           26  print(accuracy)
```



A strange curve?

As per: https://en.wikipedia.org/wiki/Talk%3ARandom_forest
(https://en.wikipedia.org/wiki/Talk%3ARandom_forest)

"Random Forests does not overfit. The testing performance of Random Forests does
not decrease (due to overfitting) as the number of trees increases. Hence after certain
number of trees the performance tend to stay in a certain value."

Microsoft

However, we can also see that ~250 estimators/trees is the ideal parameter.

Naive Bayes v Random Forest v SVM:
[https://www.researchgate.net/publication/336225950_Comparison_of_Naive_Bayes_Sup](https://www.researchgate.net/publication/336225950_Comparison_of_Naive_Bayes_Sup)
(https://www.researchgate.net/publication/336225950_Comparison_of_Naive_Bayes_Sup)

## Pull Fresh News:

In [19]:

```python
1  import requests
2  import time
3  import datetime
4
5
6  articleCount = 0
7
8  headers = {
9      'User-Agent': 'Mozilla/5.0 (Windows NT 6.1) AppleWebKit/537.36
10  }
11
12  stocks = ['TSLA', 'AMZN', 'MMM', 'INTC', 'GOOGL', 'FB', 'MSFT', 'AA
13  list_of_headlines = []
14  for line in stocks:
15      ticker = line
16
17      try:
18
19          #Query for the stock name, for refined news queries.
20          resp = requests.get(
21              url="https://www.alphavantage.co/query?function=SYMBOL_
22                  ticker), headers=headers)
23          data = resp.json()
24          companyName = data['bestMatches'][0]['2. name']
25          print("Company Name: " + companyName)
26
27          #Query for news
28          resp = requests.get(
29              url='https://newsapi.org/v2/everything?'
30  'q={}&'
31  'from=2020-01-05' # This is the OLDEST date an article can be from,
32  'sortBy=popularity&' #Filter by popularity (read the newsapi docs)
33  'apiKey=fe00115ceffe418988616191b03e1c74'.format(
34                  ticker + " " + companyName), headers=headers) #Add
35          data = resp.json()
36
37          for article in data['articles']:
38              articleCount = articleCount + 1
39              newsTitle = article['title']
40              print(newsTitle)
41              list_of_headlines.append(newsTitle)
42
43          time.sleep(1)
44
45      except Exception as e:
46          print("Error: " + str(e))
47          time.sleep(10)
48
49  # Create the pandas DataFrame and save to csv
50  df = pd.DataFrame({'headlines':list_of_headlines})
51  df.to_csv('fresh_news_month_tsla.csv', encoding='utf-8', mode='w',
```

```
Cramer Weighs In On Cracker Barrel, UPS And More
AWS Announces General Availability of Amazon Keyspaces (for Apache C
assandra)
Company Name: 3M Company
Dow Jones 378-Point Intraday Gain Fades, But 3M A Bright Spot; Netfl
ix, Tesla Weigh On Nasdaq - Investor's Business Daily
```

1x, Tesla Weigh On Nasdaq - Investor's Business Daily

Dow Jones, US Stocks Rise As Countries Begin To Reopen Economies - Investor's Business Daily

3M's stock surges on earnings beat, that was nearly 20 years in the making

3M Co (MMM) Q1 2020 Earnings Call Transcript

Is 3M Oversold At $155?

3M Company (MMM) CEO Mike Roman on Q1 2020 Results - Earnings Call Transcript

Were Hedge Funds Right About 3M Company (MMM)?

Did You Acquire 3M (MMM) Before February 9, 2017? Johnson Fistel Continues its Investigation of 3M; Should Management be Held Accountable for Investors Losses?

3M (MMM) Gains But Lags Market: What You Should Know

3M Holds Good On Its Promise To Prioritize Dividend

In [20]:

```python
 1  from nltk.corpus import stopwords
 2  from sklearn.feature_extraction.text import TfidfVectorizer
 3  from sklearn.feature_extraction.text import CountVectorizer
 4
 5  # Read in fresh news
 6  fresh_news = pd.read_csv('./fresh_news_month_tsla.csv')
 7  fresh_news['headlines'].head(5)
 8
 9  # Vectorize new text data with a max of 40 words being predictors
10  vectorizer_new_data = CountVectorizer(max_features=40, min_df=9)
11  processed_features_new_data = vectorizer_new_data.fit_transform(fre
12
13  # Vectorize training text data with a max of 40 words being predict
14  vectorizer = CountVectorizer(max_features=40, min_df=9)
15  processed_features = vectorizer.fit_transform(news_pd['text']).toar
16
17  X_train, X_test, y_train, y_test = train_test_split(processed_featu
18
19  # Predict on new/fresh news after fitting on training data
20  text_classifier = RandomForestClassifier(n_estimators=650, random_s
21  text_classifier.fit(X_train, y_train)
22
23  predictions = text_classifier.predict(processed_features_new_data)
24
25  # Output our predictions
26  print(predictions)
27
28  for i in range(len(predictions)):
29      if predictions[i] == 1:
30          print("Positive: " + fresh_news['headlines'][i])
31      if predictions[i] == -1:
32          print("Negative: " + fresh_news['headlines'][i])
33
```

```
[ 0  0  0  0  1  1  0  0 -1  0  0  0  0  0  0  0  1  0  1  1  1 -1
  0  0
  0  1  0 -1  0  0  0  0  0  0  0  0  0  0  0  0  0  1  0  0  0  0
  0  1
  1  1  1  0  1  0  0  0 -1  0  0 -1  0  0  0  0  0 -1  0  0  0  0
  0  0
  0  0  0 -1  0  0  0  0  1  0 -1  0  0  0  0 -1  1  0  1  0  0  0
  0  1
  0  0  0  0  1  0  0  0  0  0  0  1  0  0  0  0  0  0  0  1 -1  0
  0  0
  0  1  0  0 -1  0  0  0  0  0  0  0 -1  1  1  0  0  0  1  0  0  1
  0  0
  0  0  1  0  0  0  0  0  1  0  0  0  0  0  1  0]
Positive: Tesla wants to reopen California factory, but local author
ities say not yet
Positive: Ford is first auto maker to warn of lower sales, but unlik
ely to be last
Negative: Market Extra: The S&P 500 just posted the most daily swing
s of 3% or greater in more than a decade—even as the stock market hi
ts a 5-week high
Positive: The force that's propelled the stock market rally will exh
aust itself this week
Positive: Tesla Confirms Shanghai Gigafactory Shutdown, But Says I
```

Microsoft

Azure
Notebooks
(/#)

Preview
(/help/preview)

My
Projects (/q00311302/projects#)

Help
(https://docs.microsoft.com/en-
us/azure/notebooks/)

t's All 'According To Plan'
Positive: Steve Grasso Says Tesla Has Defied All Laws Of Probabilit
y'
Positive: Amazon stock hits record high on hopes for a coronavirus-r
elated boom
Negative: AMC's stock soars after report Amazon held merger talks
Positive: Netflix stock surges to record high as investors bet on st
reaming during coronavirus
Negative: Market Extra: The S&P 500 just posted the most daily swing
s of 3% or greater in more than a decade—even as the stock market hi
ts a 5-week high
Positive: Dow Jones, US Stocks Rise As Countries Begin To Reopen Eco
nomies - Investor's Business Daily
Positive: Did You Acquire 3M (MMM) Before February 9, 2017? Johnson
Fistel Continues its Investigation of 3M; Should Management be Held
Accountable for Investors Losses?
Positive: 3M (MMM) Gains But Lags Market: What You Should Know
Positive: 3M Holds Good On Its Promise To Prioritize Dividend
Positive: Stocks fights for gains as earnings season revs up - Fox B
usiness
Positive: The Importance Of Reading Footnotes – Uncovering Material
Items In Filings
Negative: Dow Blue-Chip 3M Co Surges on Coronavirus Demand, but Bond
King Warns of Danger
Negative: Why I Think You Should Buy This Defence Stock Before May
Negative: Call Traders Blast These 2 Chip Stocks
Negative: 5 Tech Stocks Poised to Beat Estimates This Earnings Seaso
n - Yahoo Finance
Positive: Making Most Of Lockdowns, Facebook Gaming Launches Earlier
Than Planned - Benzinga
Negative: Alphabet Announces First Quarter 2020 Results (Alphabet)
Negative: Market Extra: The S&P 500 just posted the most daily swing
s of 3% or greater in more than a decade—even as the stock market hi
ts a 5-week high
Positive: Big Data tech CEO on the federal government's response to
coronavirus: 'A total failure of leadership'
Positive: Dropbox's first quarterly profit is a sign of the ever-cha
nging economy
Positive: Investors have $5.1 trillion hiding out in the shares of f
ive companies, which will be tested this week
Positive: Making Most Of Lockdowns, Facebook Gaming Launches Earlier
Than Planned - Benzinga
Positive: Dropbox's first quarterly profit is a sign of the ever-cha
nging economy
Positive: Investors have $5.1 trillion hiding out in the shares of f
ive companies, which will be tested this week
Negative: FB to allow employees to work remotely until year end
Positive: How Large Option Traders Are Playing Microsoft As Cloud Bu
siness Booms - Yahoo Finance
Negative: Tech Stocks' Apr 29 Earnings Lineup: NOW, FB, MSFT, FICO,
GIB - Yahoo Finance
Negative: Hedge Funds' #3 Stock Pick Debunked Naysayers
Positive: Optimism May Be Over Done In The Equity Markets
Positive: Jonathan Angrist's Cognios Can't Deliver Despite Apple, Am
azon, Microsoft Bets
Positive: Microsoft announces registered exchange offers
Positive: The Ratings Game: Why Apple investors should be worried by

Microsoft

AT&T's earnings
Positive: Big Data tech CEO on the federal government's response to
coronavirus: 'A total failure of leadership'
Positive: Investors have $5.1 trillion hiding out in the shares of f
ive companies, which will be tested this week
Positive: Gene Munster Dismisses Goldman's Apple Downgrade, Says Cup
ertino Has Long-Term Earnings Power


Save our model to disk for production deployment to Sparkbot

In [ ]:
```python
1  from joblib import dump
2  dump(text_classifier,'sentimentclassified.joblib')
```

In [ ]:
```
1
```