# NewzTrader: Autonomous Trading Agent Implementation Using Natural Language Processing Of News Headlines

William Lyon

November 18, 2012

## Abstract

Natural Language Processing techniques are used to examine financial news headlines and generate predictions of stock price movements. News headlines from the Wall Street Journal from January 1, 2009 to November 1, 2012 are collected and paired with daily S&P 500 index returns. This information is used to train a Naive Bayes classifier and generate BUY/SELL trading signals. This trading strategy is then backtested using Wall Street Journal news headlines as a predictor for movement in the price of the S&P 500 index.
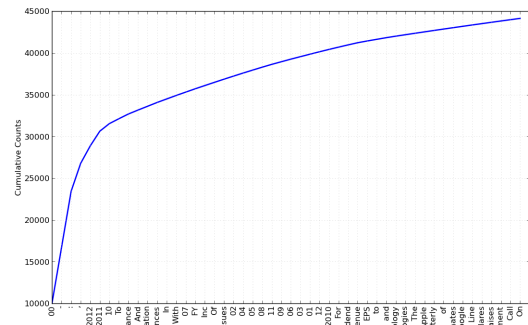
## 1 Introduction

In finance, the Efficient Market Hypothesis states that all publicly available information is reflected in financial market prices. As new information becomes available, prices adjust to take this new information into account.

### 1.1 Motivation

This tool could be used as a component of an autonomous trading agent that will make BUY/SELL decisions for trading financial instruments.

Goals: 1)

### 1.2 Literature Review

## 2 Examining The Data

Google Finance provides access to historical stock quotes as well as links to financial news stories pertaining to specific stocks. This data provides the basis for this paper.

### 2.1 Historical Stock Market Quotes

### 2.2 News Headlines

A total of 101,618 news headlines over 1416 days were collected with an average of 71 headlines per day. After aligning with financial data: UP: 14651 DOWN: 13718 NONE: 57405 TOTAL: 85775

Prior probs: UP: 0.1708 DOWN: 0.1599 NONE: 0.6693

## 3 Classification Methodology

These tuples of [return, listOfDailyNews] are used as the training data for a Naive Bayes classifier.



Figure 1: Cumulative frequency plot

### 3.1 Classes

### 3.2 Naive Bayes

### 3.3 Maximum Entropy

## 4 Implementation

### 4.1 Dependencies

### 4.2 Data Collection & Munging

### 4.3 Feature Extraction

#### 4.3.1 Bag of words

#### 4.3.2 Filtering stopwords

#### 4.3.3 Include significant bigrams

### 4.4 Naive Bayes Classifier

### 4.5 Maximum Entropy Classifier

### 4.6 Backtesting

## 5 Evaluation

### 5.1 Accuracy

### 5.2 Most informative features

### 5.3 Precision

### 5.4 Recall

## 6 Trading Model

### 6.1 Trading Signals

### 6.2 Backtesting

# 9 Source Code Listings

Listing 1: newsCredScraper.py: scrapes WSJ news headlines

```python
from lxml import etree
from datetime import datetime
from datetime import date
from datetime import timedelta
#import datetime
from dateutil.parser import parse
#import pandas as pd
import pickle
#from pandas.io.data import DataReader
#from pandas.io.data import DataReader
#from pandas.io.data import DataReader

#import pandas as pd
#from datetime import datetime

#iimport pandas as pd
#
#from pandas.io.data import DataReader

def daterange(start_date, end_date):
    for n in range(int ((end_date - start_date).days)):
        yield start_date + timedelta(n)

def only_alphanum(s):
    #s = unicode(s, "utf-8")
    return '_'.join(c for c in s.split() if c.isalnum())
def only_alpha(s):
    return '_'.join(c for c in s.split() if c.isalpha())
def removeNonAscii(s): return "".join(i for i in s if ord(i)<128)

#sp500 = DataReader("SPY", "yahoo", datetime(2009, 1, 1))
news={}
#descs = {}
start_date=date(2011, 7, 27)
end_date = date(2012, 11, 9)
news[date(2009, 11, 9)] = []
news[date(2009, 11, 10)] = []
for single_date in daterange(start_date, end_date):
    news[single_date]=[]
for single_date in daterange(start_date, end_date):
    #print "Starting date: " + str(single_date)
    year = single_date.year
    month = single_date.month
    day = single_date.day
    path = "http://api.newscred.com/articles?access_key=c4bcc3f7c9bf9ec159f51da0a86ca658&sources=104
        afa30d811d37a5582a39e1662a311&pagesize=99&from_date=%d-%d-%d&to_date=%d-%d-%d_23:59:59" % (year,
        month, day, year, month, day)
    while True:
        try:
            root = etree.parse(path)
            break
        except etree.XMLSyntaxError:
            pass
    myRoot = root.getroot()

    #news[date(year, month, day)]=[]
    #descs[date(year, month, day)]=[]
```

2

```
57        for element in myRoot.iter("article"):
58            #for item in element.iter("description"):
59              # desc = item.text
60            for item in element.iter("title"):
61                title = item.text
62            for item in element.iter("created_at"):
63                pubDate = parse(item.text)
64
65            news[pubDate.date()].append(only_alphanum(removeNonAscii(title)))
66    #import pickle
67        output = open('newsDict2.pkl', 'wb')
68        pickle.dump(news, output)
69        output.close()
```

Listing 2: dataGetter.py: downloads historical S&P 500 index price data and joins with WSJ news headlines / identifies news headline classifications and saves in format for corpus reader

```
1   # dataGetter.py
2   # William Lyon
3   # AI Grad Project
4   # NewzTrader
5   # dataGetter.py
6   # 1) Loads pickled news dict
7   # 2) Downloads historical stock price data
8   # 3) Joins news dict and stock data in a pandas dataframe
9   # 4) Set UP/DOWN/NONE classifications for every trading day
10  # 5) Save news headlines in .CSV files for corpus reader
11
12  import pickle
13  from datetime import datetime
14  from datetime import timedelta
15  from datetime import date
16
17  def daterange(start_date, end_date):
18      for n in range(int ((end_date − start_date).days)):
19          yield start_date + timedelta(n)
20
21  fkl_file = open('combinedNewsDictFull.pkl', 'rb')
22  news = pickle.load(fkl_file)
23  fkl_file.close()
24  start_date=date(2009, 1, 1)
25  end_date = date(2012, 11, 17)
26
27  import pandas as pd
28  #
29  from pandas.io.data import DataReader
30
31  def only_alphanum(s):
32      #s = unicode(s, "utf−8")
33      return ' '.join(c for c in s.split() if c.isalnum())
34  def only_alpha(s):
35      return ' '.join(c for c in s.split() if c.isalpha())
36  def removeNonAscii(s): return "".join(i for i in s if ord(i)<128)
37
38  sp500 = DataReader("^GSPC", "yahoo", datetime(2009, 1, 1))
39  newsframe = pd.Series(news, name='News')
40  #descframe = pd.Series(descs, name='Desc')
41  frameWithNews = sp500.join(pd.DataFrame(newsframe))
42  #sp500Frame = frameWithNews.join(pd.DataFrame(descframe))
43  newsframe = pd.Series(news, name='News')
44  #descframe = pd.Series(descs, name='Desc')
45  frameWithNews = sp500.join(pd.DataFrame(newsframe))
```

```
46
47   newsReturns = frameWithNews['Adj_Close'].pct_change()
48   newsReturns.name='Returns'
49   returnsFrame = frameWithNews.join(pd.DataFrame(newsReturns))
50
51   returnsFrame['UP'] = returnsFrame.Returns > 0.01
52   returnsFrame['DOWN'] = returnsFrame.Returns < −0.01
53   returnsFrame['NONE'] = (returnsFrame['UP']==False) & (returnsFrame['DOWN']==False)
54
55   droppedFrame = returnsFrame
56
57   newsUP_frame = droppedFrame[droppedFrame['UP']==True]
58
59
60   newsDOWN_frame = droppedFrame[droppedFrame['DOWN']==True]
61   newsNONE_frame = droppedFrame[droppedFrame['NONE']==True]
62
63   newsNONE_frame = newsNONE_frame.dropna()
64   newsUP_frame = newsUP_frame.dropna()
65   newsDOWN_frame = newsDOWN_frame.dropna()
66
67   # DO THIS FOR NONE, UP, and DOWN
68   i = 1
69   for row in newsNONE_frame.iterrows():
70       i+=1
71       if len(row[1].ix['News'])>0:
72           for line in row[1].ix['News']:
73               i+=1
74               writeFile = open('%d_news_NONE.csv' % i, 'w')
75               writeFile.write(line+'\n')
76
77   # DO THIS FOR NONE, UP, and DOWN
78   i = 1
79   for row in newsUP_frame.iterrows():
80       i+=1
81       if len(row[1].ix['News'])>0:
82           for line in row[1].ix['News']:
83               i+=1
84               writeFile = open('%d_news_UP.csv' % i, 'w')
85               writeFile.write(line+'\n')
86
87   # DO THIS FOR NONE, UP, and DOWN
88   i = 1
89   for row in newsDOWN_frame.iterrows():
90       i+=1
91       if len(row[1].ix['News'])>0:
92           for line in row[1].ix['News']:
93               i+=1
94               writeFile = open('%d_news_DOWN.csv' % i, 'w')
95               writeFile.write(line+'\n')
```

Listing 3: nbTrainer.py: loads news corpus / trains Naive Bayes classifier

```
1   # nbTrainer.py
2   # William Lyon
3   # AI Grad Project
4   # NewzTrader
5   # nbTrainer.py
6   # 1) Load NLTK corpus reader for naive bayes classifier
7   # 2) Train NB classifier with random 90% of corpus features
8   # 3) Test NB classifier with remaining 10% of corpus features and report
9   # accuracy
```

```python
# TODO: recall, precision reports; improve classifier (only strong words?)

# NLTK - train nb_classifier


import random
import nltk as nltk
#nltk.download()
from nltk.corpus import stopwords
import os, os.path
path = os.path.expanduser('~/nltk_data')
if not os.path.exists(path):
    os.mkdir(path)
os.path.exists(path)
import nltk.data
path in nltk.data.path
from nltk.corpus.reader import CategorizedPlaintextCorpusReader
reader = CategorizedPlaintextCorpusReader('.', r'.*_news_.*\.csv', cat_pattern=r'.*_news_(\w+)\.csv')
reader.categories()

def bag_of_words(words):
    return dict([(word, True) for word in words if word[0].isalpha()])
import collections
def bag_of_words_not_in_set(words, badwords):
    return bag_of_words(set(words)-set(badwords))

def bag_of_non_stopwords(words, stopfile='english'):
    badwords = stopwords.words(stopfile)
    return bag_of_words_not_in_set(words, badwords)

from nltk.metrics import BigramAssocMeasures
from nltk.collocations import BigramCollocationFinder

def bag_of_bigrams_words(words, score_fn=BigramAssocMeasures.chi_sq, n=2000):
    bigram_finder = BigramCollocationFinder.from_words(words)
    bigrams = bigram_finder.nbest(score_fn, n)
    dictOfBigrams = bag_of_words(bigrams)
    dictOfBigrams.update(bag_of_non_stopwords(words))
    return dictOfBigrams

def label_feats_from_corpus(corp, feature_detector=bag_of_bigrams_words):
    label_feats = collections.defaultdict(list)
    for label in corp.categories():
        for fileid in corp.fileids(categories=[label]):
            feats = feature_detector(corp.words(fileids=[fileid]))
            label_feats[label].append(feats)
    return label_feats

def split_label_feats(lfeats, split=0.90):
    train_feats = []
    test_feats = []
    for label, feats in lfeats.iteritems():
        random.shuffle(feats, random.random)
        cutoff = int(len(feats) * split)
        train_feats.extend([(feat, label) for feat in feats[:cutoff]])
        test_feats.extend([(feat, label) for feat in feats[cutoff:]])
    return train_feats, test_feats
```

```
71  reader.categories()
72
73  lfeats = label_feats_from_corpus(reader)
74  lfeats.keys()
75  train_feats, test_feats = split_label_feats(lfeats)
76  len(train_feats)
77  len(test_feats)
78
79  from nltk.classify import NaiveBayesClassifier
80  nb_classifier = NaiveBayesClassifier.train(train_feats)
81  nb_classifier.labels()
82
83  from nltk.classify.util import accuracy
84  accuracy(nb_classifier, test_feats)
```

Listing 4: backTest.py: simulates trading using trading signals generated from WSJ news headlines using Naive Bayes classifier

Listing 5: NewzTrader.py