

# Advanced Big Data - Assignment #2

Maanit Mehra - mm4694@columbia.edu

March 11, 2016

The entire code for this Assignment has been uploaded to: <https://github.com/maanitmehra/big-data>

## Problem 1

In this assignment you will be asked to predict the share price of 10 US Equity stocks at exactly 10:00 a.m. on Friday March 11.

Your predictions / code are due at exactly 8:59 a.m. March 11 (THAT MORNING)

The ticker symbols for these 10 stocks are:

BAC, C, IBM, AAPL, GE, T, MCD, NKE, TWTR, TSLA

This is an open-ended assignment so you can use as many or as few data-sources / methods as you would like in your analysis. The point of this is to do historical regression analysis to build a model and predict the stocks exactly 1 hr in advance (after 30 minutes of activity since markets open at 9:30 am). Although this may sound arbitrary the point is to implement your own creative approach to analyzing different data sources and generating tangible results. Everyone has the same tools and (random luck) for this assignment so consider it a free for all to create the best approach.

## Solution

We explore three approaches to the problem:

1. Polynomial Regression using only historical Data
2. Supervised Learning and the Least squares algorithm
3. Exploring Twitter based sentiment Analysis on Stock Behaviour.

## Part 1

The first part of the approach involves a simplified extrapolation problem. We broke this up in two stages. We first began with a basic linear model adapted over the course of 10 am collected available using the google finance API for a period of 15 days.

```
1 GOOGLE_QUOTE_URL = "http://www.google.com/finance/getprices?f=d,o,h,l,c,v&df=
  cpct&i="+str(INTERVAL*60)+"&q="
2
3 stock_url = GOOGLE_QUOTE_URL + symbol
4 resp = requests.get(stock_url)
5 resp_filter = resp.text.split(" ")
```

The shortcomings of the model revolved on the lack of sufficient historical data for calculating generic trends.

To overcome some of these limitations we explored the usage of an adaptive polynomial thresholding. We approach this by doing basic error calculations on a number of data points for a number of stocks.

The thought process behind this entails that each degree of the polynomial best describes the characteristic of the curve.

We explore the approach to the same in our analysis. More volatile stocks (Tesla, for example) are better represented by the higher order polynomial expressions.

```
1 def extrapol(x, y, x_new, degree):
2     f = np.poly1d(np.polyfit(x, y, degree))
3     y_new = f(x_new)
4     return y_new
```

This approach allows for an optimized version of polynomial regression. And lead to some interesting results.

## Part 2

During the course of this project we explored the usage of a Supervised Learning algorithm. The input parameters we consider in this scenario include Volume of trade on the previous day, the day's Highs & Lows, opening & closing prices as features.

This data can be obtained with the help of the pandas data io library.

These dataframes return Open, High, Low, Close, Volume & Adjusted Closing.

```
1 df = pio.get_data_yahoo(company, start = dt.datetime(2015, 1, 21), end = dt.
  datetime(2016, 3, 11))
2
3 npMat = df.as_matrix()
4
5 x = npMat[:-1, 1:]
```

```

6 y = npMat[1:,0]
7
8 # Using the classical Least Squares Regression Approach
9 # for quick linear regression.
10 xtx = Dot(x.T,x)
11 xy = Dot(x.T,y)
12 try:
13     w = Dot(np.linalg.inv(xtx), xy)
14 except:
15     w = Dot(np.linalg.pinv(xtx), xy)
16
17 print w

```

We observe weights such as [7.82346495e-02 -3.03513473e-02 4.75929400e-01 -1.01093090e-07] with the negative weights indicating an inverse correlation to the highs & volumes of the previous day. Potentially this implies, higher the volume of trade, less likely is the price likely to increase.

Further using cross-validation to train and test the data, we use the approach here to indicate suitable prediction values.

```

1 ## cross validation test Data
2 xtest = x[-1,:]
3 ytestAct = y[-1]
4 ytestPred= Dot(xtest, w)
5 ## Using input parameters taken from the
6 ## 10th of March as input parameters for
7 ## the 11th of March.
8
9 x_pred = npMat[-1,1:]
10 Fin_pred = Dot(x_pred,w)
11 print "%s\tPred:%.2f\t"%(company, Fin_pred)

```

### Part 3

Lastly we attempted, unfruitful in the end, the usage of Tweet & news headline based sentiment analysis.

```

1 from alchemyapi import AlchemyAPI
2
3 KEY = open('../api.key.txt', 'rb').read()[:-1]
4
5
6 def alcObj():
7     return AlchemyAPI(KEY)
8
9
10 def sentAn(obj, text):
11     response = obj.sentiment("text", text)
12     print response

```

```

13     try:
14         score = response[u'docSentiment'][u'score']
15     except:
16         score = 0
17     type = response[u'docSentiment'][u'type']
18     return score, type

```

This approach was met with two major issues.

We did not appear to find suitable historical data for a detailed tweet based approach. This limited our usage for the following. In the presence of the Twitter data, the planned approach was to involve a weighted mean of the tweet sentiment scores to the regression model. Some of our attempts at this were not very fruitful for lack of carefully calibrated twitter data.

We approached an alternative by attempting to use the information from Google & Yahoo based headlines. These approaches too however were met with lack of quality data feeds. On occasion the data feeds were not synchronised correctly.

```

1  ## Using Google News
2  GOOGLE_NEWS_URL = 'https://ajax.googleapis.com/ajax/services/search/news?v
    =1.0&userip=INSERT-USER-IP&q='
3  ## Company Names Here
4  COMPANY = 'GOOG' # for example
5
6  def sentAn(obj, text):
7      response = obj.sentiment("text", text)
8      try:
9          score = response[u'docSentiment'][u'score']
10     except:
11         ## neutral sentiments return no score.
12         score = 0
13         type = response[u'docSentiment'][u'type']
14         return score, type
15
16  def collect_data():
17      url = (GOOGLE_NEWS_URL + COMPANY)
18
19      request = urllib2.Request(url, None)
20      response = urllib2.urlopen(request)
21
22      # Process the JSON string.
23      results = simplejson.load(response)
24
25      return results
26
27  def parsing(data_collected):
28      headline_list = []
29      stories = data_collected['responseData']['results']#[

```

```
titleNoFormatting']
30     for story in stories:
31         headline_list.append( story['titleNoFormatting'])
32     return headline_list
```

We also explored the usage of an efficient method involving moving averages & flight towards the golden ratio. Codework for the same exists already, for example at: <https://github.com/lspbate/2>

We chose to not take this approach at this point.

Future explorations would however look at exploring this approach, such as <http://www.scientificpapers.org>

We hope to explore these at a later stage.

Predictions:

```
Predictions:
BAC:      13.29
C:        41.43
IBM:      140.43
AAPL:     101.08
GE:       29.89
T:        38.30
MCD:     120.06
NKE:      59.32
TWTR:     16.54
TSLA:    205.24
```