

Fall 2021

Project 1

Python Programming

Name: Ravleen Kaur

Class ID: 43

School ID: 1113138

Course: Statistics for Data Science

Course ID: DTSC-620-W01

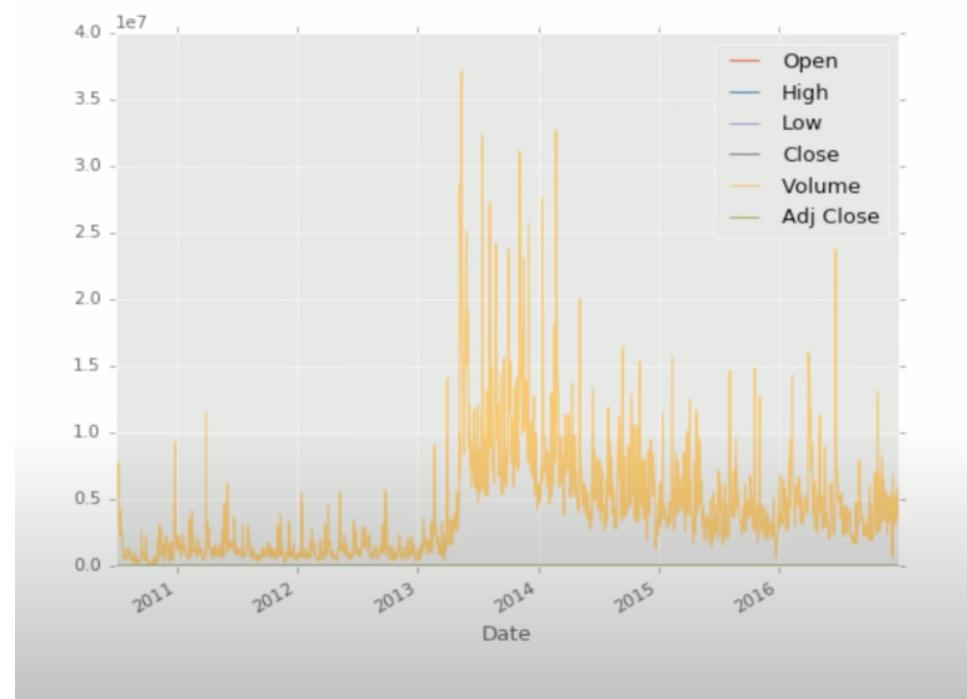
## Part 1: basic data handling with the Pandas module and data visualization with Matplotlib

Open	High	Low	Close	Volume	\
208.220001	209.990005	206.500000	208.449997	3106900	
208.000000	213.449997	207.710007	213.339996	4662900	
214.880005	222.250000	214.419998	219.529999	5901400	
221.529999	223.800003	217.199997	219.740005	3766900	
218.559998	219.199997	214.119995	214.679993	4035900	
216.300003	217.500000	211.679993	213.690002	4632700	

```
1 import datetime as dt
2 import matplotlib.pyplot as plt
3 from matplotlib import style
4 import pandas as pd
5 import pandas_datareader.data as web
6
7 style.use('ggplot')
8
9 start = dt.datetime(2000,1,1)
10 end = dt.datetime(2016,12,31)
11
12 df = web.DataReader('TSLA','yahoo',start,end)
13
14 print(df.head())
```

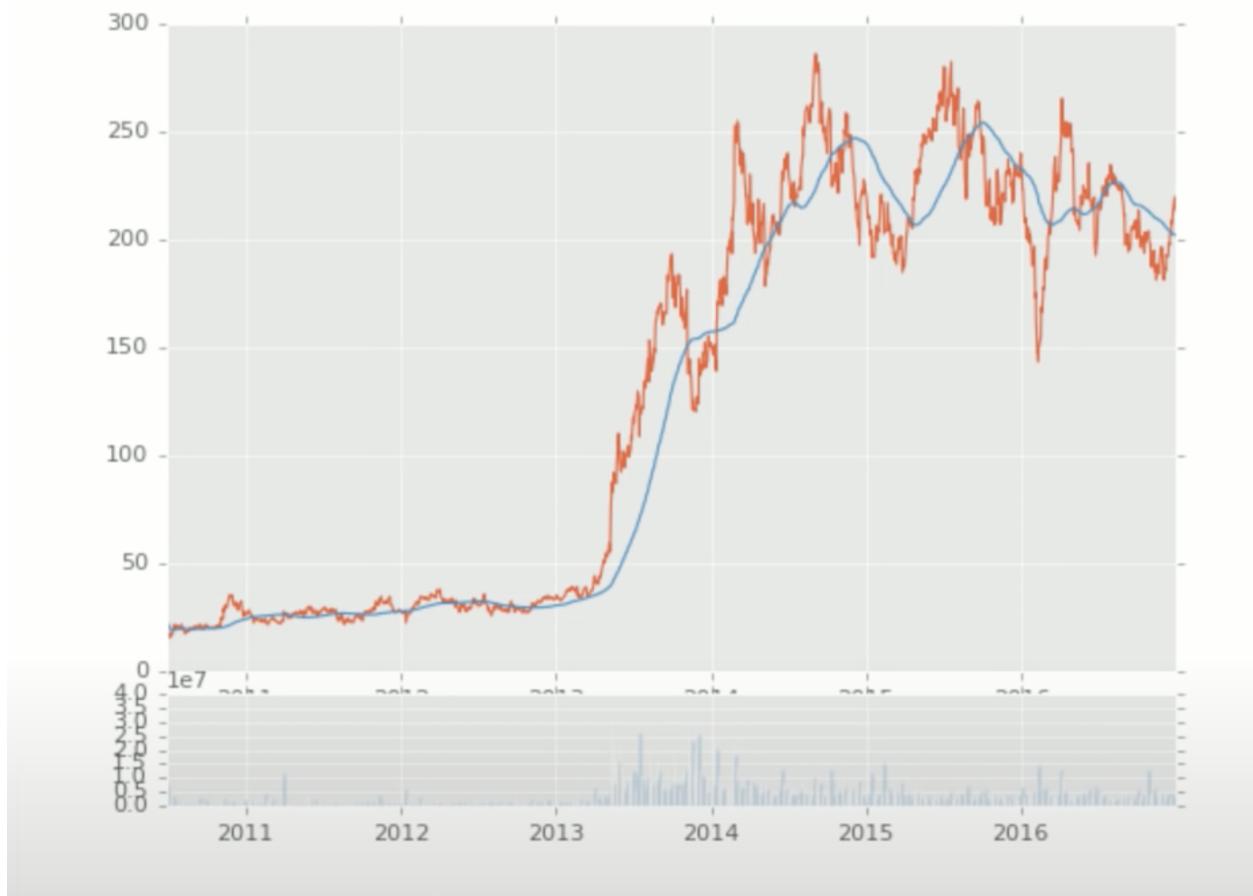
## Part 2: More data handling with the Pandas module and data visualization with Matplotlib

```
1 import datetime as dt
2 import matplotlib.pyplot as plt
3 from matplotlib import style
4 import pandas as pd
5 import pandas_datareader.data as web
6
7 style.use('ggplot')
8
9 #start = dt.datetime(2000, 1, 1)
10 #end = dt.datetime(2016, 12, 31)
11
12 #df = web.DataReader('TSLA', 'yahoo', start, end)
13 #print df.head()
14 #df.to_csv('tsla.csv')
15 df = pd.read_csv('tsla.csv', parse_dates=True, index_col=0)
16
17
18 print(df[['Open','High']].head())
19 df['Adj Close'].plot()
20 plt.show()
```



### Part 3: basic data manipulation and visualizations with our stock data

```
1 import datetime as dt
2 import matplotlib.pyplot as plt
3 from matplotlib import style
4 import pandas as pd
5 import pandas_datareader.data as web
6
7 style.use('ggplot')
8
9 #start = dt.datetime(2000, 1, 1)
10 #end = dt.datetime(2016, 12, 31)
11
12 #df = web.DataReader('TSLA', 'yahoo', start, end)
13 #print df.head()
14 #df.to_csv('tsla.csv')
15 df = pd.read_csv('tsla.csv', parse_dates=True, index_col=0)
16
17 print(df[['Open','High']].head())
18 df['100ma'] = df['Adj Close'].rolling(window=100).mean()
19
20 df.dropna(inplace=True)
21
22 print df.head()
23
24 ax1 = plt.subplot2grid((6,1),(0,0),rowspan=5, colspan=1)
25 ax2 = plt.subplot2grid((6,1),(5,0),rowspan=1, colspan=1, sharex
    =ax1)
26
27 ax1.plot(df.index, df['Adj Close'])
28 ax1.plot(df.index, df['100ma'])
29 ax2.bar (df.index, df['Volume'])
30
31 plt.show()
```



#### Part 4:

```
1 import datetime as dt
2 import matplotlib.pyplot as plt
3 from matplotlib import style
4 from matplotlib.finance import candlestick_ohlc
5 import matplotlib.dates as mdates
6 import pandas as pd
7 import pandas_datareader.data as web
8
9 style.use('ggplot')
10
11 df = pd.read_csv('tsla.csv', parse_dates=True, index_col=0)
12 #df['100ma'] = df['Adj Close'].rolling(window=100).mean()
13 #df.dropna(inplace=True)
14
15
16 df_ohlc = df['Adj Close'].resample('10D').ohlc()
17 df_volume = df['Volume'].resample('10D').sum()
18
19 df_ohlc.reset_index(inplace=True)
20
21 df_ohlc['Date'] = df_ohlc['Date'].map(mdates.date2num)
22
23 ax1 = plt.subplot2grid((6,1),(0,0), rowspan=5, colspan=1)
24 ax2 = plt.subplot2grid((6,1),(5,0), rowspan=1, colspan=1, sharex
    =ax1)
25 ax1.xaxis_date()
~~
```



```
23 ax1 = plt.subplot2grid((6,1),(0,0),rowspan=5, colspan=1)
24 ax2 = plt.subplot2grid((6,1),(5,0),rowspan=1, colspan=1, sharex
    =ax1)
25 ax1.xaxis_date()
26
27 candlestick_ohlc(ax1, df_ohlc.values, width=2, colorup='g')
28 ax2.fill_between(df_volume.index.map(mdates.date2num),df_volume
    .values, 0)
29 plt.show()
```

Part 5: grabbing pricing information en masse for a larger list of companies

```
□→ ['MMM\n', 'ABT\n', 'ABBV\n', 'ABMD\n', 'ACN\n', 'ATVI\n', 'ADBE\n', 'AMD\n', 'AAP\n', 'AES\n', 'APL\n', 'A\n', 'APD\n', 'AKAM\n', 'ALK\n', 'ALB\n', 'ARE\n', 'ALGN\n', 'ALLE\n', 'LNT\n', 'ALL\n', 'GOOGL\n', 'GOOG\n', 'MO\n', 'AMZN\n', 'AMCR\n', 'AEE\n', 'AAL\n', 'AEP\n', 'AXP\n', 'ATG\n']
```

```
1 import bs4 as bs
2 import pickle
3 import requests
4
5 def save_sp500_tickers():
6     resp = requests.get('https://en.wikipedia.org/wiki
7         /List_of_S%26P_500_companies')
8     soup = bs.BeautifulSoup(resp.text)
9     table = soup.find('table', {'class':'wikitable sortable'})
10    tickers = []
11    for row in table.findAll('tr')[1:]:
12        ticker = row.findAll('td')[0].text
13        tickers.append(str(ticker))
14    with open("sp500tickers.pickle", "wb") as f:
15        pickle.dump(tickers, f)
16
17    print tickers
18    return tickers
19
20 save_sp500_tickers()
```

```
→    'EMN\n',
    'ETN\n',
    'EBAY\n',
    'ECL\n',
    'EIX\n',
    'EW\n',
    'EA\n',
    'EMR\n',
    'ENPH\n',
    'ETR\n',
    'EOG\n',
    'EPAM\n',
    'EFX\n',
    'EQIX\n',
    'QR\n',
    'ESS\n',
    'EL\n',
    'ETSY\n',
    'EVRG\n',
    'ES\n',
    'RE\n',
    'EXC\n',
    'EXPE\n',
    'EXPD\n',
    'EXR\n',
    'XOM\n',
    'FFIV\n',
    'FAST\n',
    'FRT\n',
    'FDX\n',
    'FIS\n',
    'FITB\n',
```

Project Part 6:

```

21     print tickers
22     return tickers
23
24 #save_sp500_tickers()
25
26+ def get_data_from_yahoo(reload_sp500=False):
27+     if reload_sp500 == True:
28         tickers = save_sp500_tickers()
29+     else:
30         with open("sp500tickers.pickle", "rb") as f:
31             tickers = pickle.load(f)
32
33+     if not os.path.exists('stock_dfs'):
34         os.makedirs('stock_dfs')
35
36     start = dt.datetime(2000, 1, 1)
37     end = dt.datetime(2016, 12, 31)
38
39+     for ticker in tickers:
40         print ticker
41         if not os.path.exists('stock_dfs/{}.csv'.format(ticker)):
42             df = web.DataReader(ticker, 'yahoo', start, end)
43             df.to_csv('stock_dfs/{}.csv'.format(ticker))
44         else:
45             print('Already have {}'.format(ticker))
46
47 get_data_from_yahoo()

1 import datetime as dt
2 import os
3 import pandas as pd
4 import pandas_datareader.data as web
5 import bs4 as bs
6 import pickle
7 import requests
8
9+ def save_sp500_tickers():
10     resp = requests.get('https://en.wikipedia.org/wiki
11                         /List_of_S%26P_500_companies')
12     soup = bs.BeautifulSoup(resp.text)
13     table = soup.find('table', {'class':'wikitable sortable'})
14     tickers = []
15     for row in table.findAll('tr')[1:]:
16         ticker = row.findAll('td')[0].text
17         tickers.append(str(ticker))
18
19     with open("sp500tickers.pickle", "wb") as f:
20         pickle.dump(tickers, f)
21
22     print tickers
23     return tickers
24
25 #save_sp500_tickers()
26+ def get_data_from_yahoo(reload_sp500=False):

```

Already have ALLE  
ADS  
Already have ADS  
ALL  
Already have GOOGL  
GOOG  
Already have MO  
AMZN  
Already have AMZN  
AEE  
Already have AEE  
AAL  
Already have AAL  
AEP

Project Part 7

```
27 -     else:
28 -         with open("sp500tickers.pickle", "rb") as f:
29 -             tickers = pickle.load(f)
30
31 -     if not os.path.exists('stock_dfs'):
32 -         os.makedirs('stock_dfs')
33
34     start = dt.datetime(2000, 1, 1)
35     end = dt.datetime(2016, 12, 31)
36
37 -     for ticker in tickers:
38 -         print ticker
39 -         if not os.path.exists('stock_dfs/{}.csv'.format(ticker)):
40 -             df = web.DataReader(ticker, 'yahoo', start, end)
41 -             df.to_csv('stock_dfs/{}.csv'.format(ticker))
42 -         else:
43 -             print('Already have {}'.format(ticker))
44
45 - def compile_data():
46 -     with open("sp500tickers.pickle", "rb") as f:
47 -         tickers = pickle.load(f)
48
49     main_df = pd.DataFrame()
50
51 -     for count, ticker in enumerate(tickers):
52 -         df = pd.read_csv('/Users/kachunfung/python/finance
      /stock_dfs/{}.csv'.format(ticker))
```

```
.....           df.rename(columns = {'Adj Close': ticker}, inplace = True)
.....           df.drop(['Open', 'High', 'Low', 'Close', 'Volume'], 1,
.....           inplace=True)

.....           if main_df.empty:
.....               main_df = df
.....           else:
.....               main_df = main_df.join(df, how='outer')

.....           if count % 10 == 0:
.....               print (count)

print (main_df.head())
main_df.to_csv('sp500_joined_closes.csv')

compile_data()
```

```
import datetime as dt
import os
import pandas as pd
import pandas_datareader.data as web
import bs4 as bs
import pickle
import requests

def save_sp500_tickers():
    resp = requests.get('https://en.wikipedia.org/wiki
        /List_of_S%26P_500_companies')
    soup = bs.BeautifulSoup(resp.text)
    table = soup.find('table', {'class':'wikitable sortable'})
    tickers = []
    for row in table.findAll('tr')[1:]:
        ticker = row.findAll('td')[0].text
        tickers.append(str(ticker))

    with open("sp500tickers.pickle", "wb") as f:
        pickle.dump(tickers, f)

    print tickers
    return tickers

def get_data_from_yahoo(reload_sp500=False):
    if reload_sp500 == True:
        tickers = save_sp500_tickers()
```

```
0
10
20
30
40
50
60
70
80
90
100
110
120
130
140
150
160
170
180
190
200
210
220
230
240
250
260
270
```

## Project Part 8

```
def get_data_from_yahoo(reload_sp500=False):
    if reload_sp500 == True:
        tickers = save_sp500_tickers()
    else:
        with open("sp500tickers.pickle", "rb") as f:
            tickers = pickle.load(f)

    if not os.path.exists('stock_dfs'):
        os.makedirs('stock_dfs')

    start = dt.datetime(2000, 1, 1)
    end = dt.datetime(2016, 12, 31)

    for ticker in tickers:
        print(ticker)
        if not os.path.exists('stock_dfs/{}.csv'.format(ticker)):
            df = web.DataReader(ticker, 'yahoo', start, end)
            df.to_csv('stock_dfs/{}.csv'.format(ticker))
        else:
            print('Already have {}'.format(ticker))

def compile_data():
    with open("sp500tickers.pickle", "rb") as f:
        tickers = pickle.load(f)

import datetime as dt
import os
import pandas as pd
import pandas_datareader.data as web
import bs4 as bs
import pickle
import requests
import matplotlib.pyplot as plt
from matplotlib import style
import numpy as np

style.use('ggplot')

def save_sp500_tickers():
    resp = requests.get('https://en.wikipedia.org/wiki'
                        '/List_of_S%26P_500_companies')
    soup = bs.BeautifulSoup(resp.text)
    table = soup.find('table', {'class': 'wikitable sortable'})
    tickers = []
    for row in table.findAll('tr')[1:]:
        ticker = row.findAll('td')[0].text
        tickers.append(str(ticker))

    with open("sp500tickers.pickle", "wb") as f:
        pickle.dump(tickers, f)

    print(tickers)
```

```

# want a more natural, table-like display
ax.invert_yaxis() # top of matplotlib graph
ax.xaxis.tick_top() # move xaxis tick to the top

# set the labels
column_labels = df_corr.columns
row_labels = df_corr.index

ax.set_xticklabels(column_labels)
ax.set_yticklabels(row_labels)
plt.xticks(rotation=90) # rotating the label
heatmap.set_clim(-1, 1) # limit of color
plt.tight_layout() # clean things out, show better
plt.show()

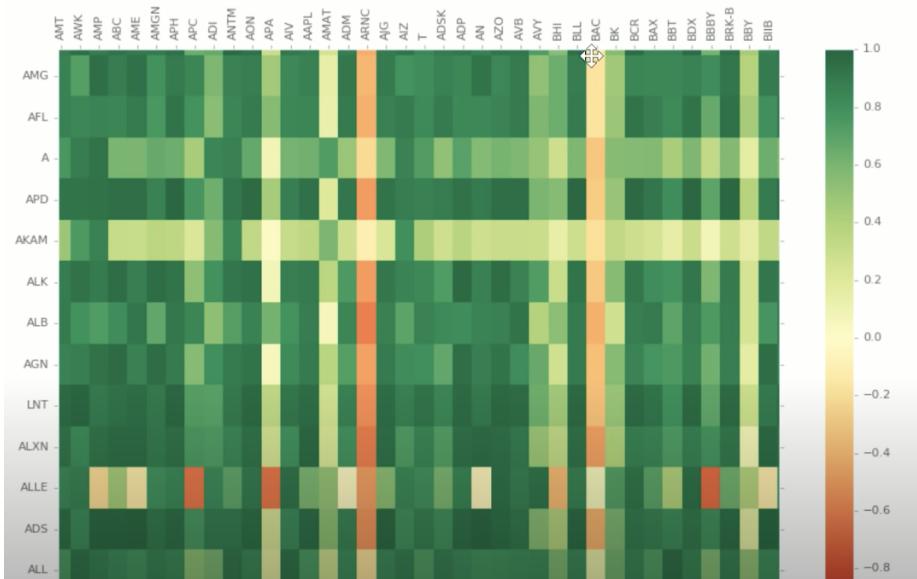
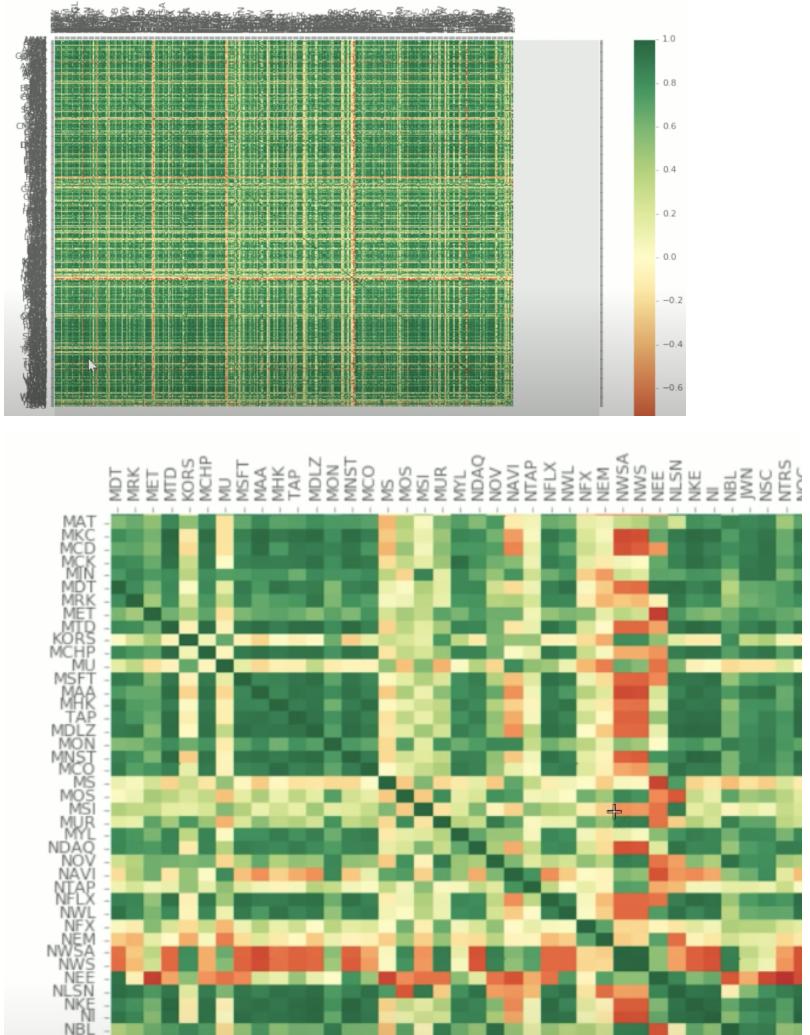
def visualize_data():
    df = pd.read_csv('sp500_joined_closes.csv')
    #df['AAPL'].plot()
    #plt.show()
    df_corr = df.corr()

    print (df_corr)
    data = df_corr.values # inner value of the dataframe, no
        index, no header
    fig = plt.figure() # create a figure object
    ax = fig.add_subplot(1,1,1) # create an axes object in the
        figure,
    # number of rows, columns, and the ID of the subplot, betw
        1 and the number of columns times the number of rows.

    heatmap = ax.pcolor(data, cmap=plt.cm.RdYlGn) #plot colour
        cmap = from Red to Yello to Green
    fig.colorbar(heatmap) # legend

    # put the major ticks at the middle of each cell
    ax.set_xticks(np.arange(data.shape[0]) + 0.5, minor=False)
    # set_xticks = Set the x ticks with list of ticks
    # arange = form an array
    # shape[0] = how many rows, shape[1] = how many columns
    # data.shape[0] = 505 rows

```



```
import numpy as np
import pandas as pd
import pickle

# 7days, 2% up buy, 2% down sell

def process_data_for_labels(ticker):
    hm_days = 7 # how many days
    df = pd.read_csv('sp500_joined_closes.csv', index_col=0)
    tickers = df.columns.values.tolist()
    df.fillna(0, inplace=True)

    for i in range(1, hm_days+1):
        # up = -i = i days in future
        # (new - old) / old
        df['{}_{0}d'.format(ticker, i)] = (df[ticker].shift(-i)
                                             / df[ticker]) - 1

    df.fillna(0, inplace=True)
    return tickers, df

print process_data_for_labels('XOM')
```

```
if col < -requirement:  
    return -1 # sell  
return 0 # hold
```

```
import numpy as np  
import pandas as pd  
import pickle  
  
# 7days, 2% up buy, 2% down sell  
  
def process_data_for_labels(ticker):  
    hm_days = 7 # how many days  
    df = pd.read_csv('sp500_joined_closes.csv', index_col=0)  
    tickers = df.columns.values.tolist()  
    df.fillna(0, inplace=True)  
  
    for i in range(1, hm_days+1):  
        # '{}'.format = String formatting  
        # up = -i = i days in future  
        # (new - old) / old  
        df['{}_{0}d'.format(ticker, i)] = (df[ticker].shift(-i) -  
            df[ticker]) / df[ticker]  
  
    df.fillna(0, inplace=True)  
    return tickers, df  
  
def buy_sell_hold(*args): # *args = any number of arguments  
    cols = [c for c in args] # list comprehension  
    requirement = 0.02 # if stock price changes for 2%  
    for col in cols:  
        if col > requirement:
```

```

plt.xticks(rotation=90)
heatmap1.set_clim(-1,1)
plt.tight_layout()
#plt.savefig("correlations.png", dpi = (300))
plt.show()

def process_data_for_labels(ticker):
    hm_days = 7
    df = pd.read_csv('sp500_joined_closes.csv', index_col=0)
    tickers = df.columns.values.tolist()
    df.fillna(0, inplace=True)

    for i in range(1,hm_days+1):
        df['{}_{ }d'.format(ticker,i)] = (df[ticker].shift(-i) - df[ticker]) / df[ticker]

    df.fillna(0, inplace=True)
    return tickers, df

def buy_sell_hold(*args):
    cols = [c for c in args]
    requirement = 0.02
    for col in cols:
        if col > requirement:
            return 1
        if col < -requirement:

def compile_data():
    with open("sp500tickers.pickle","rb") as f:
        tickers = pickle.load(f)

    main_df = pd.DataFrame()

    for count,ticker in enumerate(tickers):
        df = pd.read_csv('stock_dfs/{}.csv'.format(ticker))
        df.set_index('Date', inplace=True)

        df.rename(columns={'Adj Close':ticker}, inplace=True)
        df.drop(['Open','High','Low','Close','Volume'],1,inplace=True)

        if main_df.empty:
            main_df = df
        else:
            main_df = main_df.join(df, how='outer')

        if count % 10 == 0:
            print(count)
    print(main_df.head())
    main_df.to_csv('sp500_joined_closes.csv')

```

```

        (ticker)],
        df['{}_5d'.format
(ticker)],
        df['{}_6d'.format
(ticker)],
        df['{}_7d'.format
(ticker)] ))
```

vals = df['{}\_target'.format(ticker)].values.tolist()  
str\_vals = [str(i) for i in vals]  
print('Data spread:', str(Counter(str\_vals)))

df.fillna(0, inplace=True)  
df = df.replace([np.inf, -np.inf], np.nan)  
df.dropna(inplace=True)

df\_vals = df[[ticker\_name for ticker\_name in tickers]]  
.pct\_change()  
df\_vals = df\_vals.replace([np.inf, -np.inf], 0)  
df\_vals.fillna(0, inplace=True)

X = df\_vals.values  
y = df['{}\_target'.format(ticker)].values  
return X,y,df

```
print extract_featuresets('XOM')
```

```
Data spread: Counter({'1': 1713, '-1': 1456, '0': 1108})
```

```
import bs4 as bs
from collections import Counter
import datetime as dt
import matplotlib.pyplot as plt
from matplotlib import style
import numpy as np
import os
import pandas as pd
import pandas_datareader.data as web
import pickle
import requests
from sklearn import svm, cross_validation, neighbors
from sklearn.ensemble import VotingClassifier, RandomForestClassifier
import statistics

style.use('ggplot')

def save_sp500_tickers():
    resp = requests.get('http://en.wikipedia.org/wiki/List_of_S%26P_500_companies')
    soup = bs.BeautifulSoup(resp.text, 'lxml')
    table = soup.find('table', {'class': 'wikitable sortable'})
    tickers = []
    for row in table.findAll('tr')[1:]:
        ticker = row.findAll('td')[0].text
        tickers.append(ticker)

    with open("sp500tickers.pickle","wb") as f:
        pickle.dump(tickers,f)

    return tickers
```

```
def get_data_from_yahoo(reload_sp500=False):

    if reload_sp500:
        tickers = save_sp500_tickers()
    else:
        with open("sp500tickers.pickle","rb") as f:
            tickers = pickle.load(f)

    if not os.path.exists('stock_dfs'):
        os.makedirs('stock_dfs')

    start = dt.datetime(2000, 1, 1)
    end = dt.datetime(2016, 12, 31)

    for ticker in tickers:
        # just in case your connection breaks, we'd like to save our progress!
        if not os.path.exists('stock_dfs/{}.csv'.format(ticker)):
            df = web.DataReader(ticker, "yahoo", start, end)
            df.to_csv('stock_dfs/{}.csv'.format(ticker))
        else:
            print('Already have {}'.format(ticker))

def compile_data():
    with open("sp500tickers.pickle","rb") as f:
        tickers = pickle.load(f)

    main_df = pd.DataFrame()

    for count,ticker in enumerate(tickers):
        df = pd.read_csv('stock_dfs/{}.csv'.format(ticker))
        df.set_index('Date', inplace=True)
```

```
        main_df = df
    else:
        main_df = main_df.join(df, how='outer')

    if count % 10 == 0:
        print(count)
# print(main_df.head())
main_df.to_csv('sp500_joined_closes.csv')

def visualize_data():
    df = pd.read_csv('sp500_joined_closes.csv')
#df['AAPL'].plot()
#plt.show()
    df_corr = df.corr()
# print(df_corr.head())
    df_corr.to_csv('sp500corr.csv')

    data1 = df_corr.values
    fig1 = plt.figure()
    ax1 = fig1.add_subplot(111)

    heatmap1 = ax1.pcolor(data1, cmap=plt.cm.RdYlGn)
    fig1.colorbar(heatmap1)

    ax1.set_xticks(np.arange(data1.shape[1]) + 0.5, minor=False)
    ax1.set_yticks(np.arange(data1.shape[0]) + 0.5, minor=False)
    ax1.invert_xaxis()
```

```

def extract_featuresets(ticker):
    tickers, df = process_data_for_labels(ticker)

    df['{}_target'.format(ticker)] = list(map( buy_sell_hold,
                                              df['{}_1d'.format(ticker)],
                                              df['{}_2d'.format(ticker)],
                                              df['{}_3d'.format(ticker)],
                                              df['{}_4d'.format(ticker)],
                                              df['{}_5d'.format(ticker)],
                                              df['{}_6d'.format(ticker)],
                                              df['{}_7d'.format(ticker)] ))

    vals = df['{}_target'.format(ticker)].values.tolist()
    str_vals = [str(i) for i in vals]
    print('Data spread:',Counter(str_vals))

    df.fillna(0, inplace=True)
    df = df.replace([np.inf, -np.inf], np.nan)
    df.dropna(inplace=True)

    df_vals = df[[ticker_name for ticker_name in tickers]].pct_change()
    df_vals = df_vals.replace([np.inf, -np.inf], 0)
    df_vals.fillna(0, inplace=True)

    X = df_vals.values
    y = df['{}_target'.format(ticker)].values

    return X,y,df

```

```

#clf = neighbors.KNeighborsClassifier()

clf = VotingClassifier([('lsvc',svm.LinearSVC()),
                       ('knn',neighbors.KNeighborsClassifier()),
                       ('rfor',RandomForestClassifier())])

clf.fit(X_train, y_train)
confidence = clf.score(X_test, y_test)
print('accuracy:',confidence)
predictions = clf.predict(X_test)
print('predicted class counts:',Counter(predictions))
print()
print()
return confidence

from statistics import mean

with open("sp500tickers.pickle","rb") as f:
    tickers = pickle.load(f)

accuracies = []
for count,ticker in enumerate(tickers):

    if count%10==0:
        print(count)

    accuracy = do_ml(ticker)
    accuracies.append(accuracy)
    print("{} accuracy: {}. Average accuracy:{}".format(ticker,accuracy,mean(accuracies)))

```