

# New York Stock Price Prediction with RNN

This project demonstrates the future price prediction for different stocks using Recurrent Neural Networks in Tensorflow.

```
In [4]: import numpy as np
import pandas as pd
import math
import sklearn
import sklearn.preprocessing
import datetime
import os
import matplotlib.pyplot as plt
import tensorflow as tf
```

```
In [ ]: #display parent directory and working directory
print(os.path.dirname(os.getcwd())+':', os.listdir(os.path.dirname(os.getcwd())));
print(os.getcwd()+':', os.listdir(os.getcwd()));

# split data in 80%/10%/10% train/validation/test sets
valid_set_size_percentage = 10
test_set_size_percentage = 10
```

Analyze data

```
In [6]: # import all stock prices
df = pd.read_csv("../input/prices-split-adjusted.csv", index_col = 0)
df.info()

# number of different stocks
print('\nnumber of different stocks: ', len(list(set(df.symbol))))
print(list(set(df.symbol))[:10])
```

```

<class 'pandas.core.frame.DataFrame'>
Index: 851264 entries, 2016-01-05 to 2016-12-30
Data columns (total 6 columns):
symbol      851264 non-null object
open        851264 non-null float64
close       851264 non-null float64
low         851264 non-null float64
high        851264 non-null float64
volume      851264 non-null float64
dtypes: float64(5), object(1)
memory usage: 45.5+ MB

number of different stocks: 501
['MLM', 'ABC', 'TXN', 'EL', 'NTAP', 'PNW', 'VRTX', 'IBM', 'OKE', 'GPN']

```

In [7]: `df.tail()`

```

Out[7]:

```

	symbol	open	close	low	high	volume
<b>date</b>						
<b>2016-12-30</b>	ZBH	103.309998	103.199997	102.849998	103.930000	973800.0
<b>2016-12-30</b>	ZION	43.070000	43.040001	42.689999	43.310001	1938100.0
<b>2016-12-30</b>	ZTS	53.639999	53.529999	53.270000	53.740002	1701200.0
<b>2016-12-30</b>	AIV	44.730000	45.450001	44.410000	45.590000	1380900.0
<b>2016-12-30</b>	FTV	54.200001	53.630001	53.389999	54.480000	705100.0

In [8]: `df.describe()`

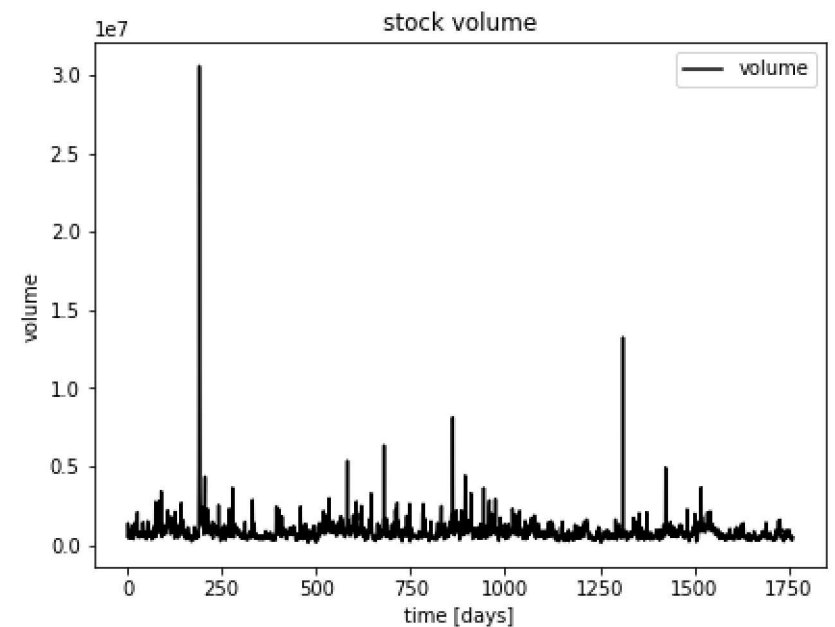
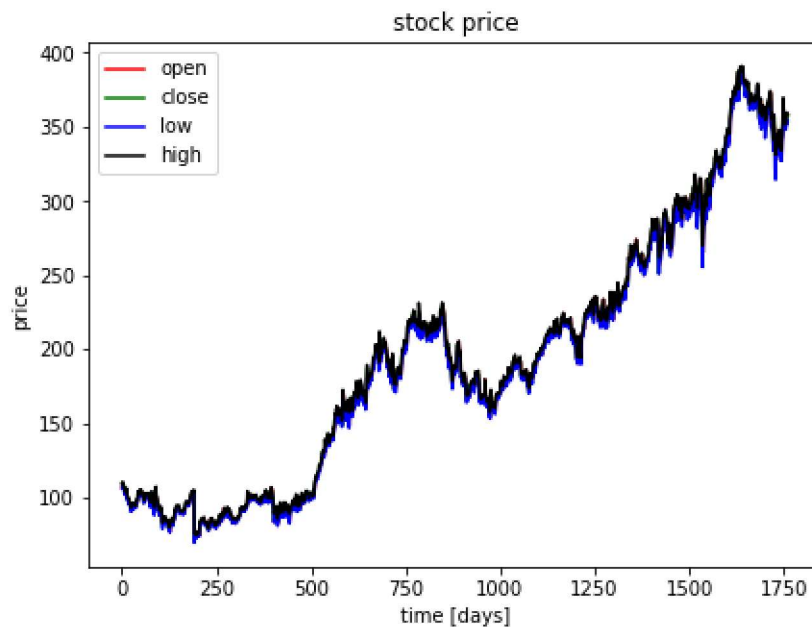
Out[8]:

	open	close	low	high	volume
<b>count</b>	851264.000000	851264.000000	851264.000000	851264.000000	8.512640e+05
<b>mean</b>	64.993618	65.011913	64.336541	65.639748	5.415113e+06
<b>std</b>	75.203893	75.201216	74.459518	75.906861	1.249468e+07
<b>min</b>	1.660000	1.590000	1.500000	1.810000	0.000000e+00
<b>25%</b>	31.270000	31.292776	30.940001	31.620001	1.221500e+06
<b>50%</b>	48.459999	48.480000	47.970001	48.959999	2.476250e+06
<b>75%</b>	75.120003	75.139999	74.400002	75.849998	5.222500e+06
<b>max</b>	1584.439941	1578.130005	1549.939941	1600.930054	8.596434e+08

In [10]:

```
plt.figure(figsize=(15, 5));
plt.subplot(1,2,1);
plt.plot(df[df.symbol == 'EQIX'].open.values, color='red', label='open')
plt.plot(df[df.symbol == 'EQIX'].close.values, color='green', label='close')
plt.plot(df[df.symbol == 'EQIX'].low.values, color='blue', label='low')
plt.plot(df[df.symbol == 'EQIX'].high.values, color='black', label='high')
plt.title('stock price')
plt.xlabel('time [days]')
plt.ylabel('price')
plt.legend(loc='best')
#plt.show()

plt.subplot(1,2,2);
plt.plot(df[df.symbol == 'EQIX'].volume.values, color='black', label='volume')
plt.title('stock volume')
plt.xlabel('time [days]')
plt.ylabel('volume')
plt.legend(loc='best');
```



Manipulate data

- choose a specific stock
- drop feature: volume
- normalize stock data
- create train, validation and test data sets

```
In [11]: # function for min-max normalization of stock
def normalize_data(df):
    min_max_scaler = sklearn.preprocessing.MinMaxScaler()
    df['open'] = min_max_scaler.fit_transform(df.open.values.reshape(-1,1))
    df['high'] = min_max_scaler.fit_transform(df.high.values.reshape(-1,1))
    df['low'] = min_max_scaler.fit_transform(df.low.values.reshape(-1,1))
    df['close'] = min_max_scaler.fit_transform(df['close'].values.reshape(-1,1))
    return df

# function to create train, validation, test data given stock data and sequence length
def load_data(stock, seq_len):
    data_raw = stock.as_matrix() # convert to numpy array
    data = []

    # create all possible sequences of length seq_len
```

```

for index in range(len(data_raw) - seq_len):
    data.append(data_raw[index: index + seq_len])

data = np.array(data);
valid_set_size = int(np.round(valid_set_size_percentage/100*data.shape[0]));
test_set_size = int(np.round(test_set_size_percentage/100*data.shape[0]));
train_set_size = data.shape[0] - (valid_set_size + test_set_size);

x_train = data[:train_set_size,:-1,:];
y_train = data[:train_set_size,-1,:];

x_valid = data[train_set_size:train_set_size+valid_set_size,:-1,:];
y_valid = data[train_set_size:train_set_size+valid_set_size,-1,:];

x_test = data[train_set_size+valid_set_size:,:-1,:];
y_test = data[train_set_size+valid_set_size:,-1,:];

return [x_train, y_train, x_valid, y_valid, x_test, y_test]

# choose one stock
df_stock = df[df.symbol == 'EQIX'].copy()
df_stock.drop(['symbol'],1,inplace=True)
df_stock.drop(['volume'],1,inplace=True)

cols = list(df_stock.columns.values)
print('df_stock.columns.values = ', cols)

# normalize stock
df_stock_norm = df_stock.copy()
df_stock_norm = normalize_data(df_stock_norm)

# create train, test data
seq_len = 20 # choose sequence length
x_train, y_train, x_valid, y_valid, x_test, y_test = load_data(df_stock_norm, seq_len)
print('x_train.shape = ',x_train.shape)
print('y_train.shape = ', y_train.shape)
print('x_valid.shape = ',x_valid.shape)
print('y_valid.shape = ', y_valid.shape)
print('x_test.shape = ', x_test.shape)
print('y_test.shape = ',y_test.shape)

```

```

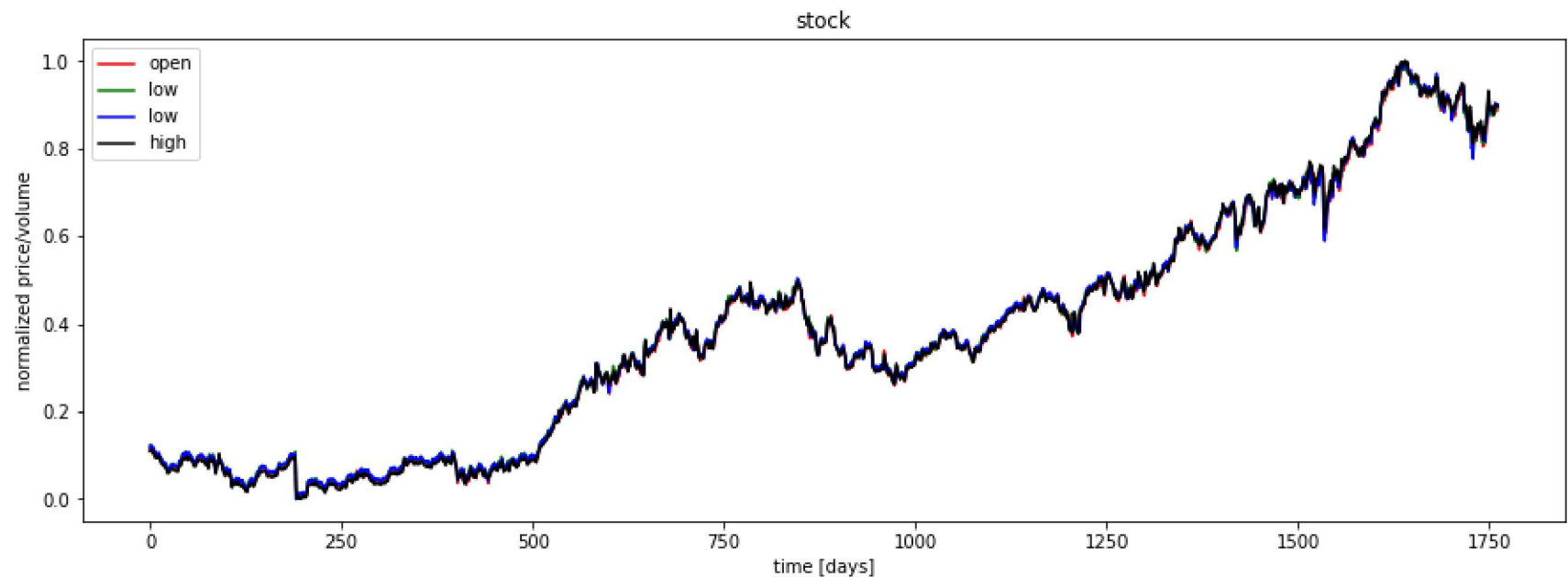
df_stock.columns.values = ['open', 'close', 'low', 'high']
x_train.shape = (1394, 19, 4)
y_train.shape = (1394, 4)
x_valid.shape = (174, 19, 4)
y_valid.shape = (174, 4)
x_test.shape = (174, 19, 4)
y_test.shape = (174, 4)

```

```

In [12]: plt.figure(figsize=(15, 5));
plt.plot(df_stock_norm.open.values, color='red', label='open')
plt.plot(df_stock_norm.close.values, color='green', label='low')
plt.plot(df_stock_norm.low.values, color='blue', label='low')
plt.plot(df_stock_norm.high.values, color='black', label='high')
#plt.plot(df_stock_norm.volume.values, color='gray', label='volume')
plt.title('stock')
plt.xlabel('time [days]')
plt.ylabel('normalized price/volume')
plt.legend(loc='best')
plt.show()

```



Model and validate data

- RNN with basic

In [13]: *## Basic RNN in Tensorflow*

```
index_in_epoch = 0;
perm_array = np.arange(x_train.shape[0])
np.random.shuffle(perm_array)

# function to get the next batch
def get_next_batch(batch_size):
    global index_in_epoch, x_train, perm_array
    start = index_in_epoch
    index_in_epoch += batch_size

    if index_in_epoch > x_train.shape[0]:
        np.random.shuffle(perm_array) # shuffle permutation array
        start = 0 # start next epoch
        index_in_epoch = batch_size

    end = index_in_epoch
    return x_train[perm_array[start:end]], y_train[perm_array[start:end]]

# parameters
n_steps = seq_len-1
n_inputs = 4
n_neurons = 200
n_outputs = 4
n_layers = 2
learning_rate = 0.001
batch_size = 50
n_epochs = 100
train_set_size = x_train.shape[0]
test_set_size = x_test.shape[0]

tf.reset_default_graph()

X = tf.placeholder(tf.float32, [None, n_steps, n_inputs])
y = tf.placeholder(tf.float32, [None, n_outputs])

# use Basic RNN Cell
layers = [tf.contrib.rnn.BasicRNNCell(num_units=n_neurons, activation=tf.nn.elu)
          for layer in range(n_layers)]

# use Basic LSTM Cell
# layers = [tf.contrib.rnn.BasicLSTMCell(num_units=n_neurons, activation=tf.nn.elu)
#           for layer in range(n_layers)]
```

```

# use LSTM Cell with peephole connections
#layers = [tf.contrib.rnn.LSTMCell(num_units=n_neurons,
#                                  activation=tf.nn.leaky_relu, use_peepholes = True)
#          for layer in range(n_Layers)]

# use GRU cell
#layers = [tf.contrib.rnn.GRUCell(num_units=n_neurons, activation=tf.nn.leaky_relu)
#          for layer in range(n_Layers)]

multi_layer_cell = tf.contrib.rnn.MultiRNNCell(layers)
rnn_outputs, states = tf.nn.dynamic_rnn(multi_layer_cell, X, dtype=tf.float32)

stacked_rnn_outputs = tf.reshape(rnn_outputs, [-1, n_neurons])
stacked_outputs = tf.layers.dense(stacked_rnn_outputs, n_outputs)
outputs = tf.reshape(stacked_outputs, [-1, n_steps, n_outputs])
outputs = outputs[:,n_steps-1,:] # keep only last output of sequence

loss = tf.reduce_mean(tf.square(outputs - y)) # loss function = mean squared error
optimizer = tf.train.AdamOptimizer(learning_rate=learning_rate)
training_op = optimizer.minimize(loss)

# run graph
with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    for iteration in range(int(n_epochs*train_set_size/batch_size)):
        x_batch, y_batch = get_next_batch(batch_size) # fetch the next training batch
        sess.run(training_op, feed_dict={X: x_batch, y: y_batch})
        if iteration % int(5*train_set_size/batch_size) == 0:
            mse_train = loss.eval(feed_dict={X: x_train, y: y_train})
            mse_valid = loss.eval(feed_dict={X: x_valid, y: y_valid})
            print('%0.2f epochs: MSE train/valid = %0.6f/%0.6f'%(
                iteration*batch_size/train_set_size, mse_train, mse_valid))

    y_train_pred = sess.run(outputs, feed_dict={X: x_train})
    y_valid_pred = sess.run(outputs, feed_dict={X: x_valid})
    y_test_pred = sess.run(outputs, feed_dict={X: x_test})

```



```
0.00 epochs: MSE train/valid = 0.390838/0.489912
4.99 epochs: MSE train/valid = 0.000193/0.000892
9.97 epochs: MSE train/valid = 0.000144/0.000780
14.96 epochs: MSE train/valid = 0.000130/0.000693
19.94 epochs: MSE train/valid = 0.000118/0.000499
24.93 epochs: MSE train/valid = 0.000100/0.000406
29.91 epochs: MSE train/valid = 0.000122/0.000335
34.90 epochs: MSE train/valid = 0.000086/0.000318
39.89 epochs: MSE train/valid = 0.000080/0.000319
44.87 epochs: MSE train/valid = 0.000080/0.000271
49.86 epochs: MSE train/valid = 0.000163/0.000673
54.84 epochs: MSE train/valid = 0.000091/0.000265
59.83 epochs: MSE train/valid = 0.000089/0.000253
64.81 epochs: MSE train/valid = 0.000077/0.000283
69.80 epochs: MSE train/valid = 0.000088/0.000437
74.78 epochs: MSE train/valid = 0.000070/0.000256
79.77 epochs: MSE train/valid = 0.000067/0.000236
84.76 epochs: MSE train/valid = 0.000068/0.000241
89.74 epochs: MSE train/valid = 0.000070/0.000245
94.73 epochs: MSE train/valid = 0.000095/0.000269
99.71 epochs: MSE train/valid = 0.000068/0.000195
```

Predictions

```
In [14]: y_train.shape
```

```
Out[14]: (1394, 4)
```

```
In [15]: ft = 0 # 0 = open, 1 = close, 2 = highest, 3 = lowest

## show predictions
plt.figure(figsize=(15, 5));
plt.subplot(1,2,1);

plt.plot(np.arange(y_train.shape[0]), y_train[:,ft], color='blue', label='train target')

plt.plot(np.arange(y_train.shape[0], y_train.shape[0]+y_valid.shape[0]), y_valid[:,ft],
         color='gray', label='valid target')

plt.plot(np.arange(y_train.shape[0]+y_valid.shape[0],
                  y_train.shape[0]+y_test.shape[0]+y_test.shape[0]),
         y_test[:,ft], color='black', label='test target')

plt.plot(np.arange(y_train_pred.shape[0]),y_train_pred[:,ft], color='red',
```

```

        label='train prediction')

plt.plot(np.arange(y_train_pred.shape[0], y_train_pred.shape[0]+y_valid_pred.shape[0]),
         y_valid_pred[:,ft], color='orange', label='valid prediction')

plt.plot(np.arange(y_train_pred.shape[0]+y_valid_pred.shape[0],
                  y_train_pred.shape[0]+y_valid_pred.shape[0]+y_test_pred.shape[0]),
         y_test_pred[:,ft], color='green', label='test prediction')

plt.title('past and future stock prices')
plt.xlabel('time [days]')
plt.ylabel('normalized price')
plt.legend(loc='best');

plt.subplot(1,2,2);

plt.plot(np.arange(y_train.shape[0], y_train.shape[0]+y_test.shape[0]),
         y_test[:,ft], color='black', label='test target')

plt.plot(np.arange(y_train_pred.shape[0], y_train_pred.shape[0]+y_test_pred.shape[0]),
         y_test_pred[:,ft], color='green', label='test prediction')

plt.title('future stock prices')
plt.xlabel('time [days]')
plt.ylabel('normalized price')
plt.legend(loc='best');

corr_price_development_train = np.sum(np.equal(np.sign(y_train[:,1]-y_train[:,0]),
        np.sign(y_train_pred[:,1]-y_train_pred[:,0])).astype(int)) / y_train.shape[0]
corr_price_development_valid = np.sum(np.equal(np.sign(y_valid[:,1]-y_valid[:,0]),
        np.sign(y_valid_pred[:,1]-y_valid_pred[:,0])).astype(int)) / y_valid.shape[0]
corr_price_development_test = np.sum(np.equal(np.sign(y_test[:,1]-y_test[:,0]),
        np.sign(y_test_pred[:,1]-y_test_pred[:,0])).astype(int)) / y_test.shape[0]

print('correct sign prediction for close - open price for train/valid/test: %.2f/%.2f/%.2f'%(
    corr_price_development_train, corr_price_development_valid, corr_price_development_test))

```

correct sign prediction for close - open price for train/valid/test: 0.75/0.64/0.61

