

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import datetime
4 import seaborn as sns
5
6 import matplotlib.pyplot as plt
7 import matplotlib.mlab as mlab
8 import matplotlib
9 plt.style.use('ggplot')
10 from matplotlib.pyplot import figure
11
12 %matplotlib inline
13 matplotlib.rcParams['figure.figsize'] = (12,8)
14
15 pd.options.mode.chained_assignment = None
16
17 from collections import Counter
```

```
In [2]: 1 # read the data
2 df = pd.read_csv('all_stocks_5yr.csv')
3
4 # shape and data types of the data
5 print("Shape of our data is: ", df.shape)
6 print()
7
8 # type of each data
9 print("Data type is: ")
10
11 dType = pd.DataFrame(df.dtypes, columns = ['Type'])
12
13 dType
```

Shape of our data is: (619040, 7)

Data type is:

Out[2]:

	Type
date	object
open	float64
high	float64
low	float64
close	float64
volume	int64
Name	object

In [3]: 1 df

Out[3]:

	date	open	high	low	close	volume	Name
0	2013-02-08	15.07	15.12	14.63	14.75	8407500	AAL
1	2013-02-11	14.89	15.01	14.26	14.46	8882000	AAL
2	2013-02-12	14.45	14.51	14.10	14.27	8126000	AAL
3	2013-02-13	14.30	14.94	14.25	14.66	10259500	AAL
4	2013-02-14	14.94	14.96	13.16	13.99	31879900	AAL
...
619035	2018-02-01	76.84	78.27	76.69	77.82	2982259	ZTS
619036	2018-02-02	77.53	78.12	76.73	76.78	2595187	ZTS
619037	2018-02-05	76.64	76.92	73.18	73.83	2962031	ZTS
619038	2018-02-06	72.74	74.56	72.13	73.27	4924323	ZTS
619039	2018-02-07	72.70	75.00	72.69	73.86	4534912	ZTS

619040 rows × 7 columns

In [4]: 1 # 1. Clean the data

In [5]: 1 # ** Missing data **

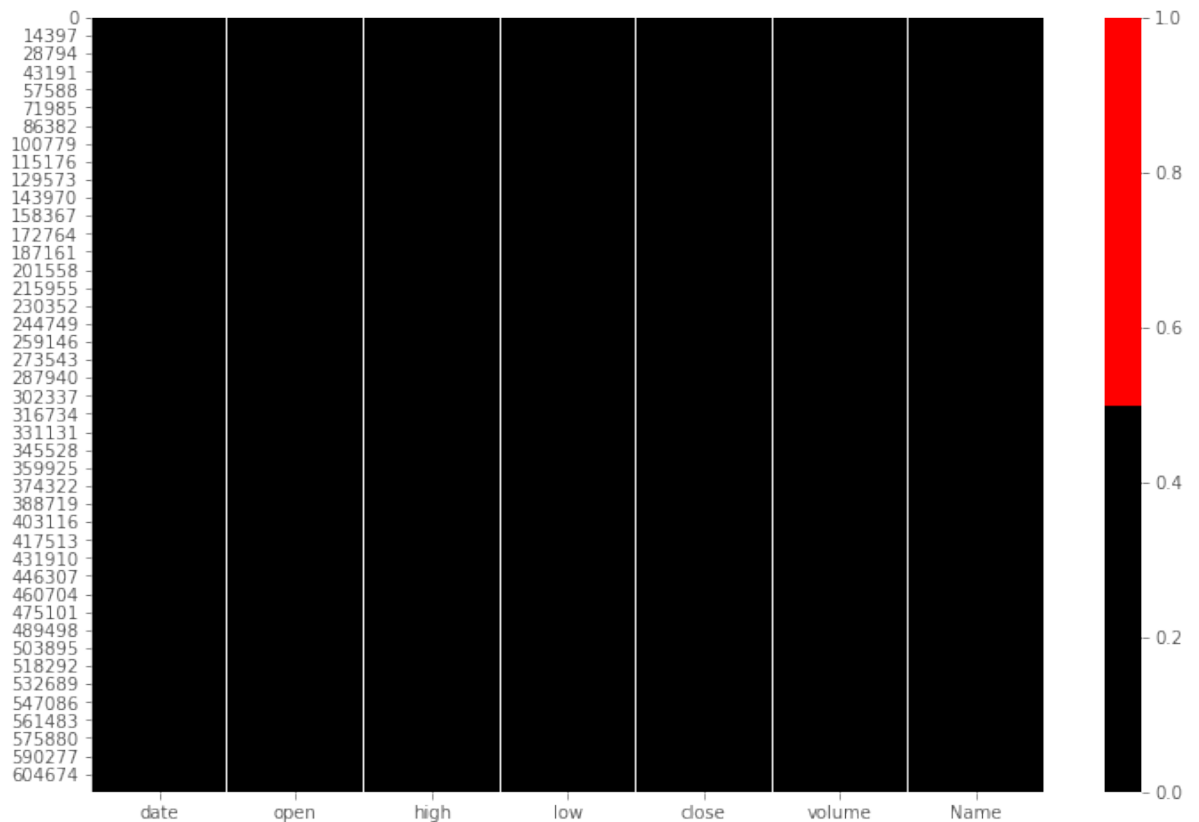
In [6]: 1 # Display Column with missing value
2 print("Number of missing data per column: ")
3 df.isnull().sum()

Number of missing data per column:

Out[6]: date 0
open 11
high 8
low 8
close 0
volume 0
Name 0
dtype: int64

```
In [7]: 1 # Visualization: number of missing value is in red, black is no
2 sns.heatmap(df[df.columns[0:7]].isnull(), cmap=sns.color_palett
3
4 # Notice that open has 11 missing value and high and low all ha
5 # value that we can't see it in the visualization below
```

Out[7]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb973b8cac0>



```
In [8]: 1 # To handle with missing data, we decide to impute the missing
2 # of the same feature that's not missing. There are 619040 diff
3 # only less than 12 missing values. It's such a small number th
4 # impute it
5 df['open'] = df['open'].fillna(df['open'].median())
6 df['high'] = df['high'].fillna(df['high'].median())
7 df['low'] = df['low'].fillna(df['low'].median())
```

```
In [9]: 1 # Display Column with missing value after handling missing valu
2 print("Number of missing data per column after handling missing
3 df.isnull().sum()
```

Number of missing data per column after handling missing values:

```
Out[9]: date      0
open      0
high      0
low       0
close     0
volume    0
Name      0
dtype: int64
```

In []:

1

In [10]:

1 *# ** Outliers ***

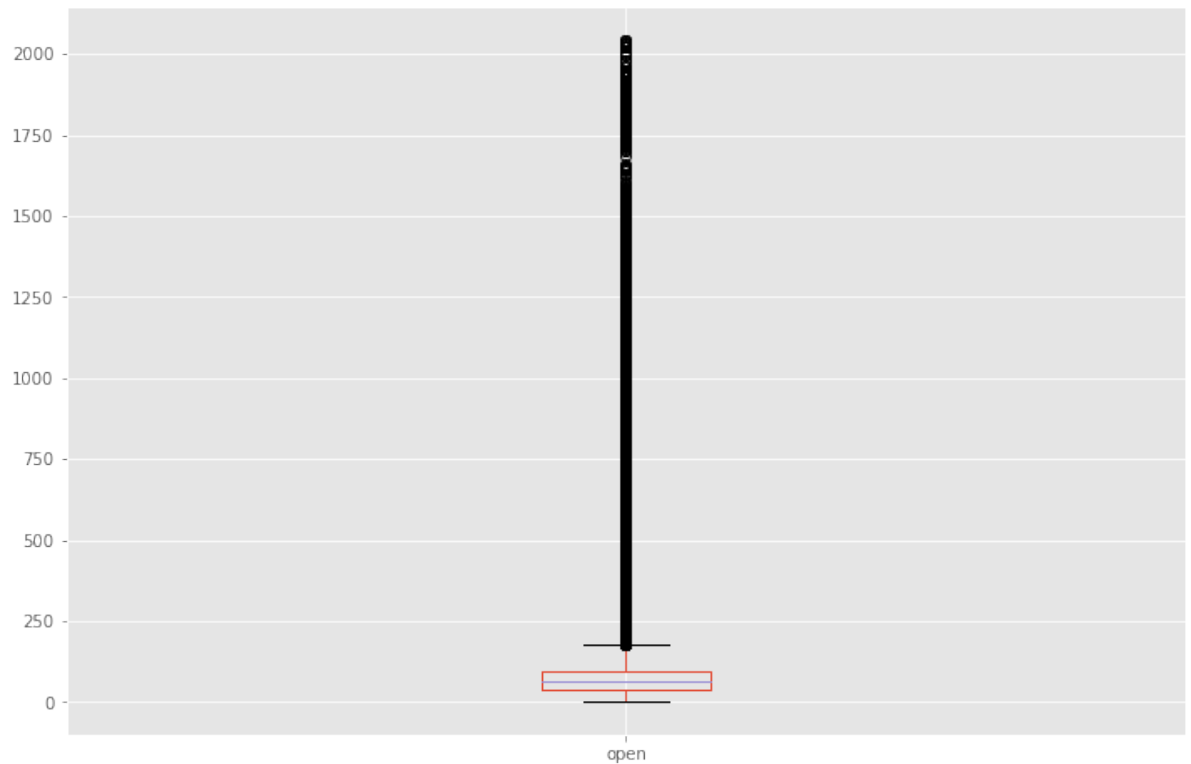
In [11]:

1 *# Observing from our data, we can see that there are 5 columns*
2 *# open, high, low, close, and volume*

In [12]:

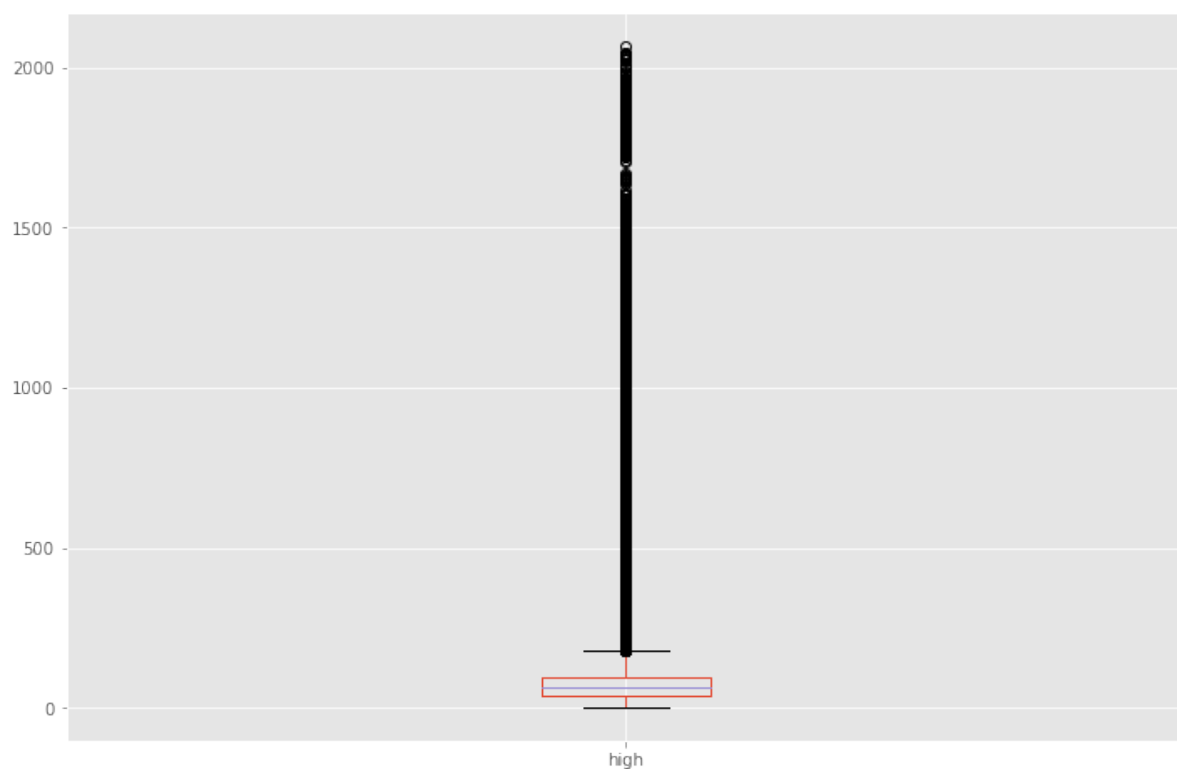
1 *# We'll start with open. There seems to be no outliers*
2 `df.boxplot(column=['open'])`

Out[12]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb959e4e5b0>



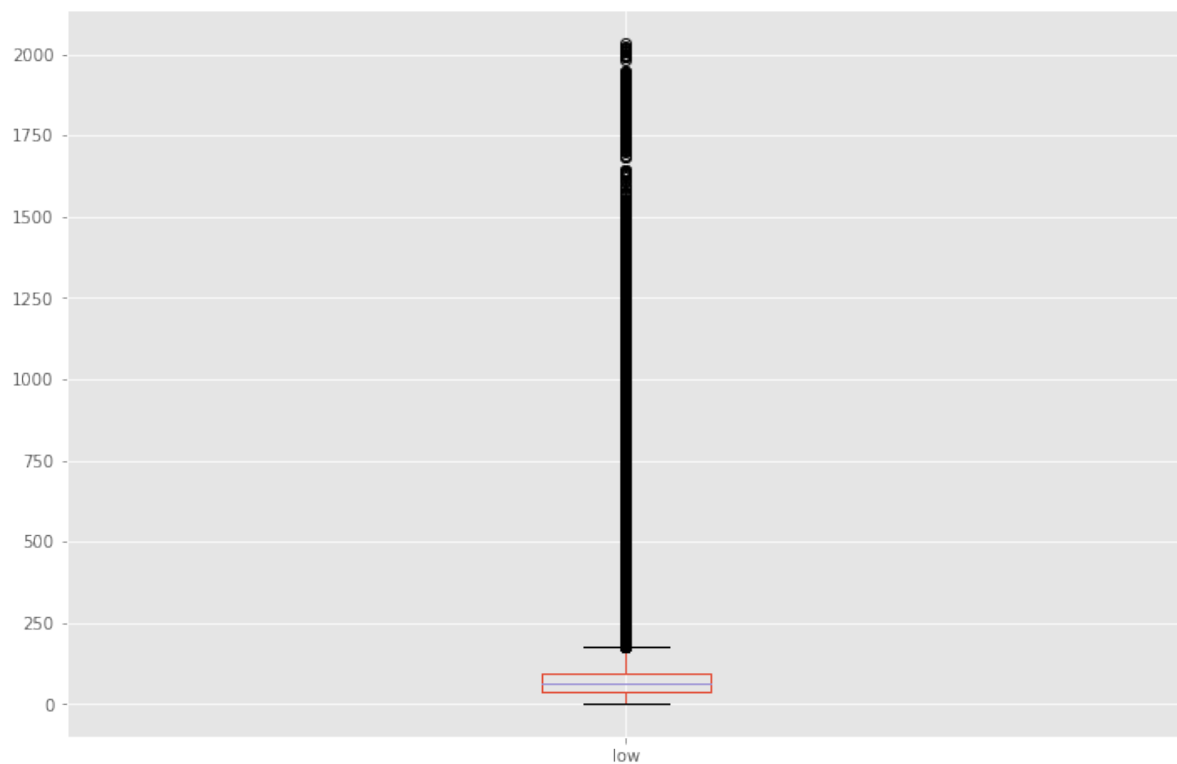
```
In [13]: 1 # Next is high. There seems to be no outliers
         2 df.boxplot(column=['high'])
```

Out[13]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb959fd77c0>



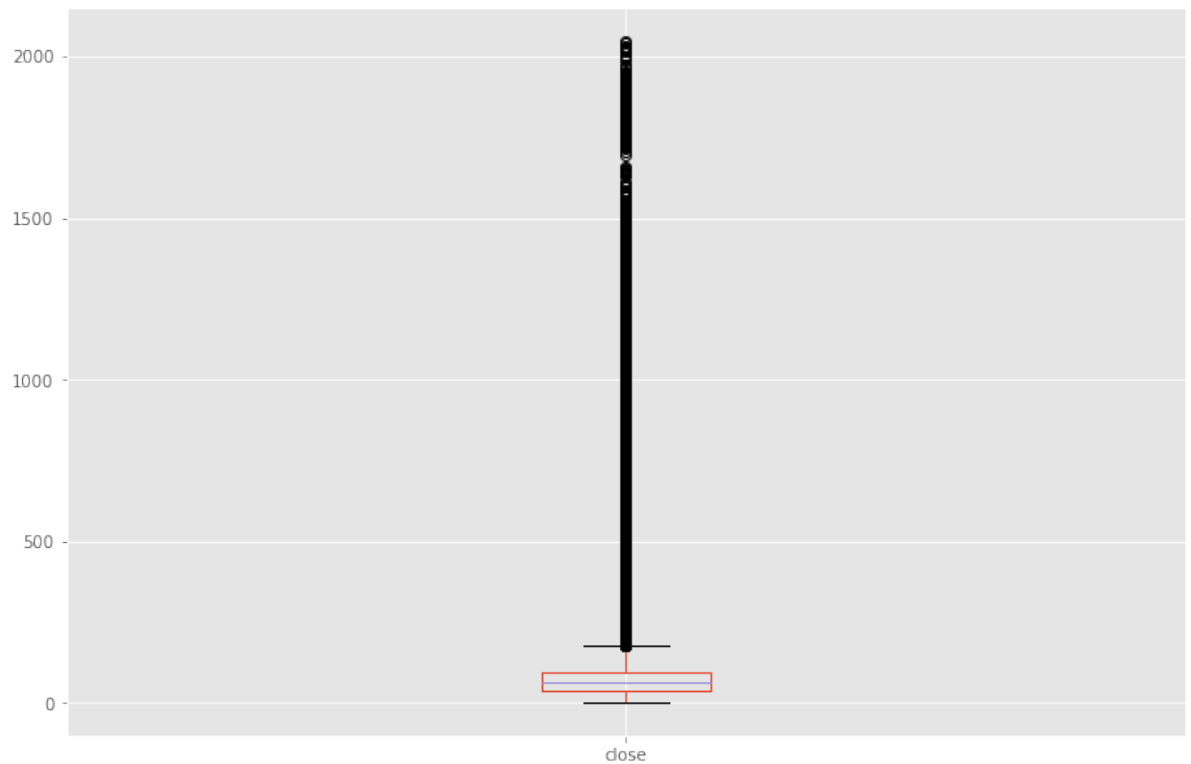
```
In [14]: 1 # Next is high. There seems to be no outliers
         2 df.boxplot(column=['low'])
```

Out[14]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb95a3b47f0>



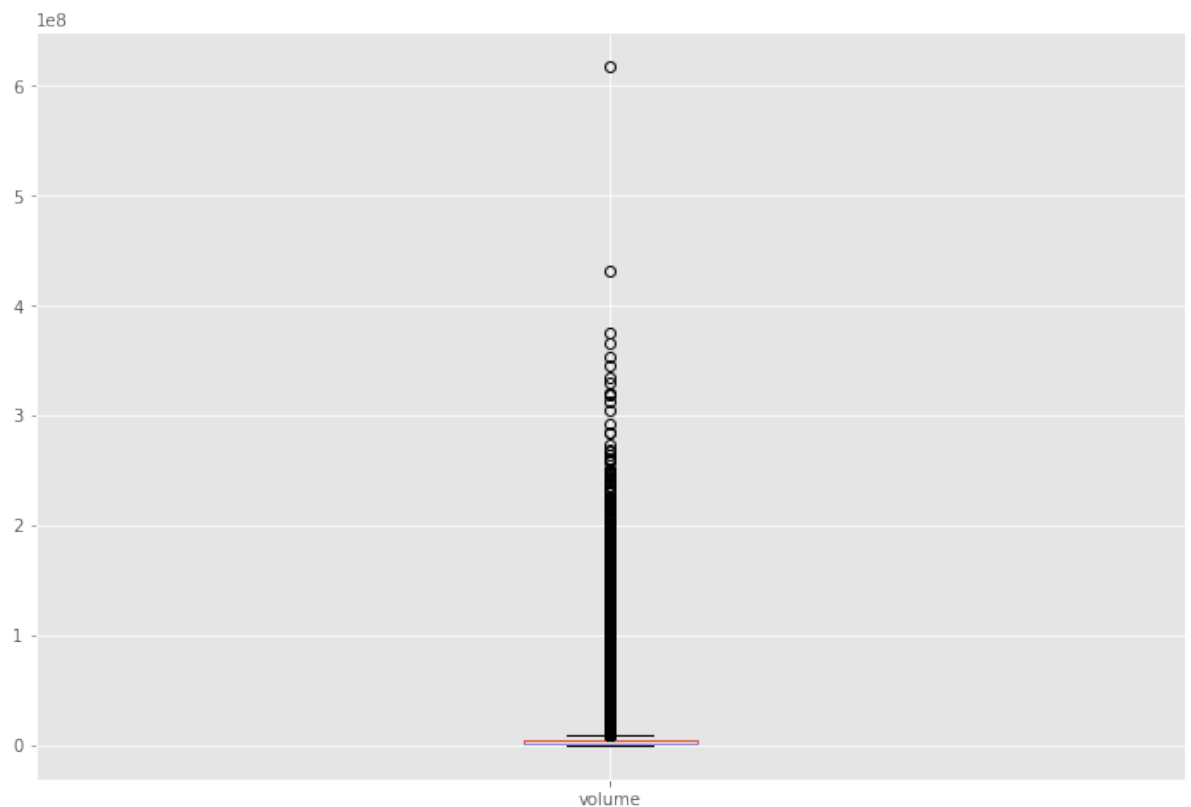
```
In [15]: 1 # Next is close. There seems to be no outliers
         2 df.boxplot(column=['close'])
```

Out[15]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb95bc54610>



```
In [16]: 1 # Next is volume.
         2 df.boxplot(column=['volume'])
```

Out[16]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb95c193b20>



```
In [17]: 1 # Seems like we have an outlier. We can look at the descriptive
        2 df['volume'].describe()
```

```
Out[17]: count      6.190400e+05
         mean      4.321823e+06
         std       8.693610e+06
         min       0.000000e+00
         25%       1.070320e+06
         50%       2.082094e+06
         75%       4.284509e+06
         max       6.182376e+08
         Name: volume, dtype: float64
```

```
In [18]: 1 # As you can see, the maximum of 6.182376e+08 is our outlier be
        2 # However, we decide to keep with outlier because it's an impor
        3 # spike in a certain day, which can indicate good news and bad
```

```
In [ ]: 1
```

```
In [19]: 1 # ** Noisy data **
        2 # First, we will check for repetitive data. I want to find colu
        3 for col in df.columns:
        4     if (df[col].value_counts(dropna=False)/len(df.index)).iloc[
        5         print('{0}: {1:.5f}%'.format(col, (df[col].value_counts
        6         print()
        7
        8 # Nothing got printed, which means more than 90% of the data ar
```

```
In [20]: 1 # Second, we will check for irrelevant value. we have skimmed t
        2 # high and low don't really provide any valuable information fo
        3 # For this project, we only interest in the open and close valu
        4 # the end of the day
        5 df = df.drop(['high', 'low'], axis=1)
        6
        7 # Now we only have 5 columns
        8 print("Shape of our data is: ", df.shape)
        9 print()
```

Shape of our data is: (619040, 5)

```
In [21]: 1 # Third, we will check for duplicated data by checking if there
        2 # date and name together shoule be unique as it's describing a
        3 if df[df.duplicated(subset=['Name', 'date'])].empty:
        4     print("There's no duplicated row")
        5 else:
        6     print("There're duplicated rows")
        7
        8 # Seems like we have no duplicated row
```

There's no duplicated row

```
In [ ]: 1
```

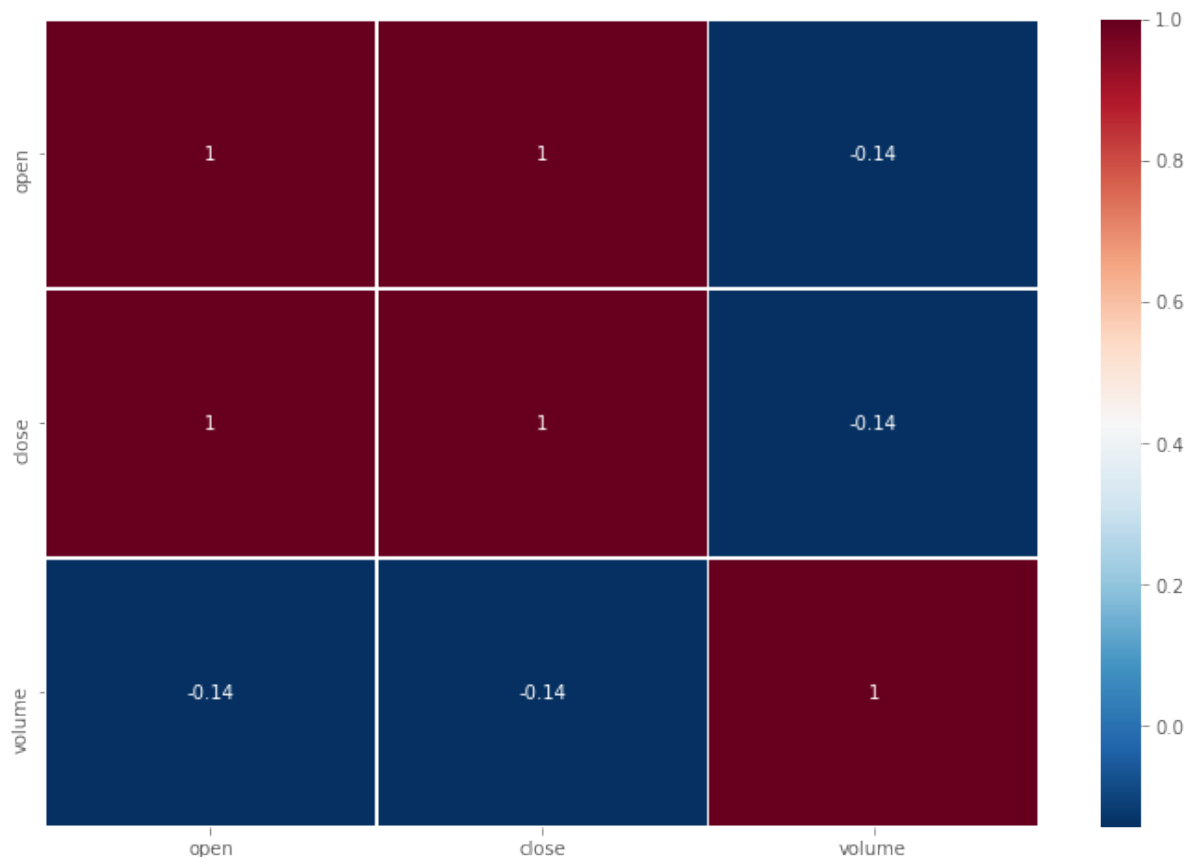
```
In [22]: 1 # ** Inconsistent Data **
```

```
In [23]: 1 # First, inconsistent usage of upper and lower cases in categor  
2 # put all letters to upper cases  
3 df['Name'] = df['Name'].str.upper()
```

```
In [24]: 1 # Second, column date doesn't have the correct data format so I  
2 # object to datetime for easier analysis later
```

```
In [25]: 1 from sklearn.linear_model import LinearRegression  
2 from sklearn.model_selection import train_test_split  
3 import matplotlib.pyplot as plt  
4 import numpy as np  
5 import seaborn as sb  
6  
7 # we check for correlation between the features using the Pears  
8 corr = df.corr(method='pearson')  
9  
10 # print the correlation table  
11 corr  
12  
13 # we print the heatmap of the correaltion table  
14 sb.heatmap(corr,xticklabels=corr.columns, yticklabels=corr.colu  
15 cmap='RdBu_r', annot=True, linewidth=0.5)  
16  
17 # The dark maroon cells represent highly correlated features of
```

```
Out[25]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb9503a7a60>
```

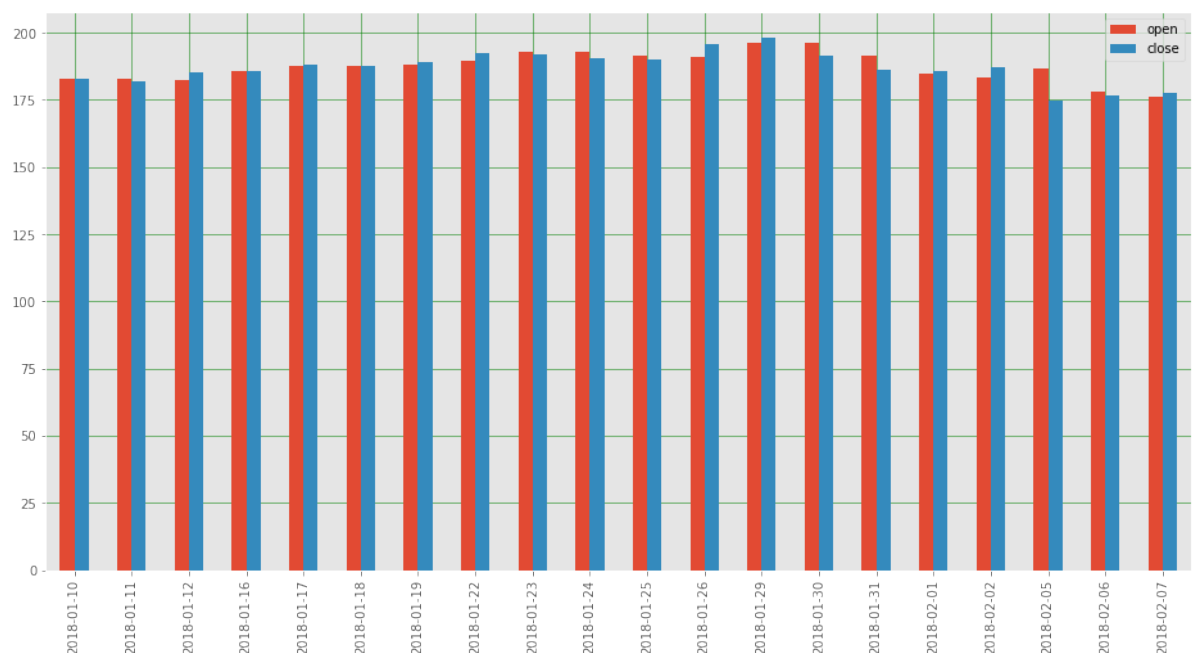



```

In [26]: 1 # We design a function to visualize the independent and depende
2
3 def OpenClosePlot(stock, name):
4     stock_df=stock[['date', 'open', 'close']]
5
6     #Plot Open vs Close
7     tickvalues = range(0,20)
8     stock_df[['open', 'close']].tail(20).plot(kind='bar',figsize
9     plt.grid(which='major', linestyle='-', linewidth='0.5', col
10    plt.grid(which='minor', linestyle=':', linewidth='0.5', col
11    plt.xticks(tickvalues, stock.date[len(stock)-20:])
12    plt.show()
13
14 Stock = input("Stock Name: ").upper()
15 stock_df = df.loc[df['Name'] == Stock]
16 OpenClosePlot(stock_df, Stock)

```

Stock Name: amgn



```

In [27]: 1 # Model training and testing
2 stock_df['date'] = pd.to_datetime(stock_df['date'], format='%Y-%
3
4 stock_df['year']=stock_df['date'].dt.year
5 stock_df['month']=stock_df['date'].dt.month
6 stock_df['day']=stock_df['date'].dt.day
7
8 stock_df = stock_df[['day', 'month', 'year', 'open', 'close']]
9 stock_df.head(10)
10
11 X = stock_df.iloc[:,stock_df.columns != 'close']
12 Y= stock_df.iloc[:, 4]
13
14 Y = Y.reset_index(drop = True)
15
16 # X = X.set_index(X['day'] + X['month'] + X['year'])
17 # X

```

```

In [28]: 1 # We train the model after testing and splitting the data
2 from sklearn.model_selection import train_test_split
3 from sklearn import model_selection
4 from sklearn.model_selection import KFold
5 import sklearn.metrics as sm
6
7 x_train,x_test,y_train,y_test= train_test_split(X,Y,test_size=.
8
9 # Use linear regression to fit the training data
10 model=LinearRegression()
11 model.fit(x_train,y_train)
12
13 # predict the values using k-fold
14 y_pred=model.predict(x_test)
15 kfold = model_selection.KFold(n_splits=5)
16 results_kfold = model_selection.cross_val_score(model, x_test,
17 print("K-Fold Accuracy: ", results_kfold.mean()*100)
18
19 print("Score Accuracy: ", model.score(x_test, y_test))
20
21 print("Mean absolute error =", round(sm.mean_absolute_error(y_t
22 print("Mean squared error =", round(sm.mean_squared_error(y_tes
23 print("Median absolute error =", round(sm.median_absolute_error
24 print("Explain variance score =", round(sm.explained_variance_s
25 print("R2 score =", round(sm.r2_score(y_test, y_pred), 2))
26
27 plot=pd.DataFrame({'Actual':y_test,'Pred':y_pred})
28 plot.head(20).plot(kind='bar',figsize=(16,8))
29 plt.grid(which='major', linestyle='-', linewidth='0.5', color='
30 plt.grid(which='minor', linestyle=':', linewidth='0.5', color='
31 plt.xlabel("Linear Predicted vs Actual")
32 plt.show()

```

K-Fold Accuracy: 99.42777021462138

Score Accuracy: 0.9944887855030116

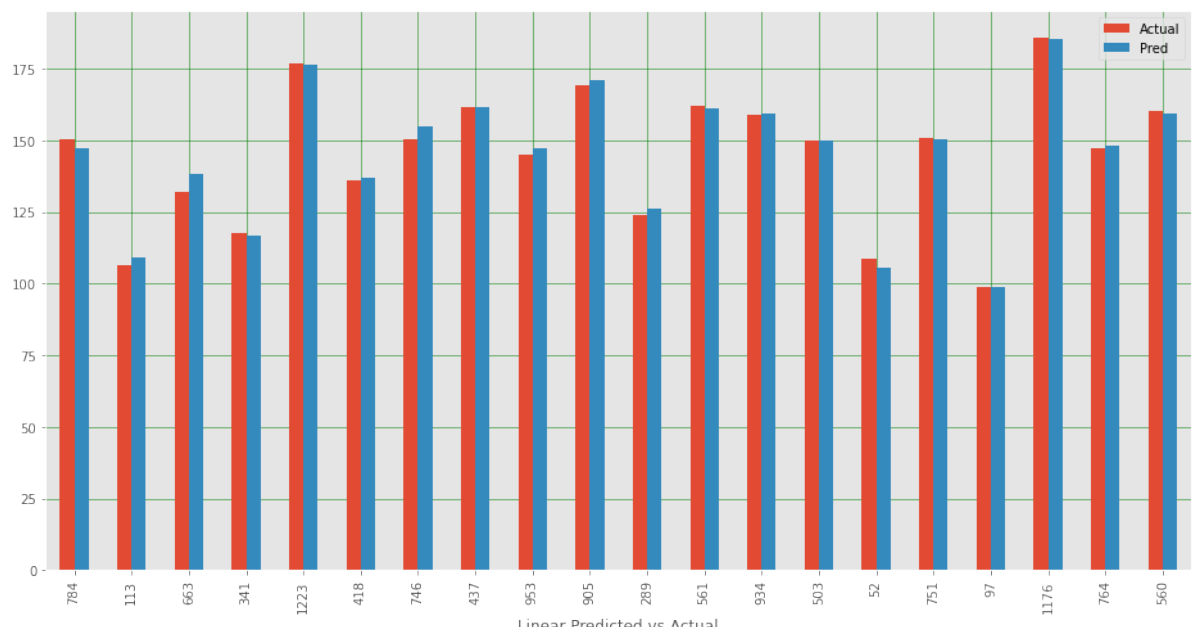
Mean absolute error = 1.41

Mean squared error = 3.71

Median absolute error = 1.03

Explain variance score = 0.99

R2 score = 0.99



In [29]:

```
1 # KNN model training
2 from sklearn.neighbors import KNeighborsRegressor
3
4 knn_regressor=KNeighborsRegressor(n_neighbors = 5)
5 knn_model=knn_regressor.fit(x_train,y_train)
6 y_knn_pred=knn_model.predict(x_test)
7
8 knn_kfold = model_selection.KFold(n_splits=20)
9 results_kfold = model_selection.cross_val_score(knn_model, x_te
10 print("K-Fold Accuracy: ", results_kfold.mean()*100)
11
12 print("Score Accuracy: ", knn_model.score(x_test, y_test))
13
14 print("Mean absolute error =", round(sm.mean_absolute_error(y_t
15 print("Mean squared error =", round(sm.mean_squared_error(y_tes
16 print("Median absolute error =", round(sm.median_absolute_error
17 print("Explain variance score =", round(sm.explained_variance_s
18 print("R2 score =", round(sm.r2_score(y_test, y_knn_pred), 2))
19
20 plot_knn_df=pd.DataFrame({'Actual':y_test,'Pred':y_knn_pred})
21 plot_knn_df.head(20).plot(kind='bar',figsize=(16,8))
22 plt.grid(which='major', linestyle='-', linewidth='0.5', color='
23 plt.grid(which='minor', linestyle=':', linewidth='0.5', color='
24 plt.xlabel("kNN Predicted vs Actual")
25 plt.show()
```

K-Fold Accuracy: 98.68568551726746

Score Accuracy: 0.9932669296352159

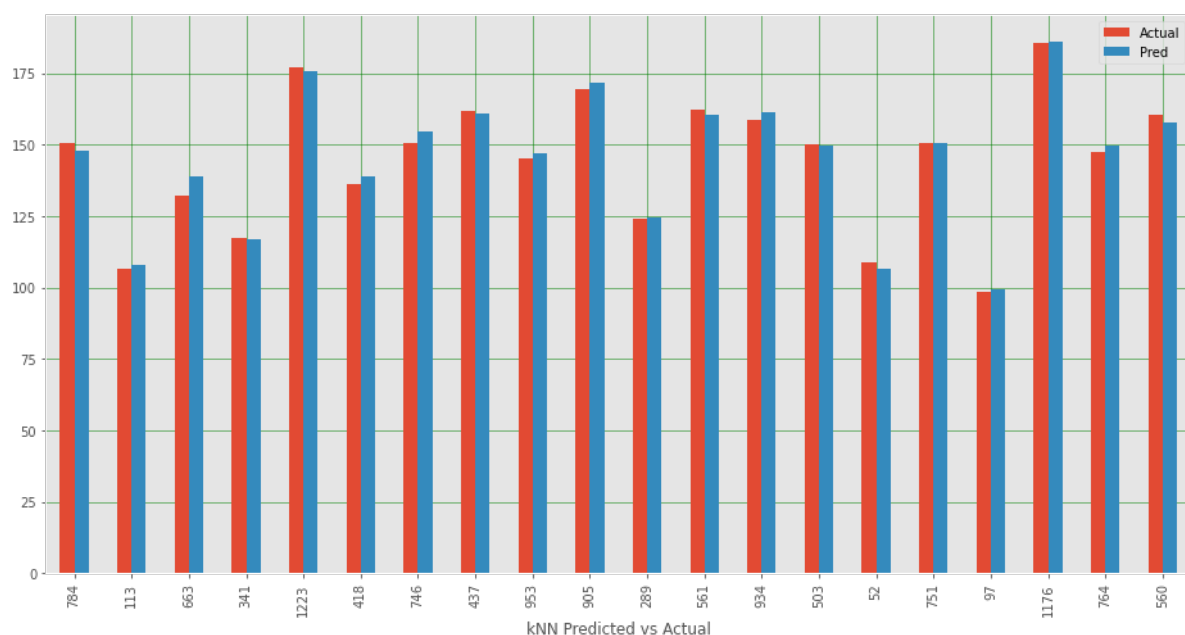
Mean absolute error = 1.63

Mean squared error = 4.54

Median absolute error = 1.33

Explain variance score = 0.99

R2 score = 0.99



In [30]:

```
1 from sklearn.svm import SVR
2 svm_regressor = SVR(kernel='linear')
3 svm_model=svm_regressor.fit(x_train,y_train)
4 y_svm_pred=svm_model.predict(x_test)
5
6 print("Score Accuracy: ", svm_model.score(x_test, y_test))
7
8 print("Mean absolute error =", round(sm.mean_absolute_error(y_t
9 print("Mean squared error =", round(sm.mean_squared_error(y_tes
10 print("Median absolute error =", round(sm.median_absolute_error
11 print("Explain variance score =", round(sm.explained_variance_s
12 print("R2 score =", round(sm.r2_score(y_test, y_svm_pred), 2))
13
14 plot_knn_df=pd.DataFrame({'Actual':y_test,'Pred':y_knn_pred})
15 plot_knn_df.head(20).plot(kind='bar',figsize=(16,8))
16 plt.grid(which='major', linestyle='-', linewidth='0.5', color='
17 plt.grid(which='minor', linestyle=':', linewidth='0.5', color='
18 plt.xlabel("SVM Predicted vs Actual")
19 plt.show()
```

Score Accuracy: 0.994264659399689

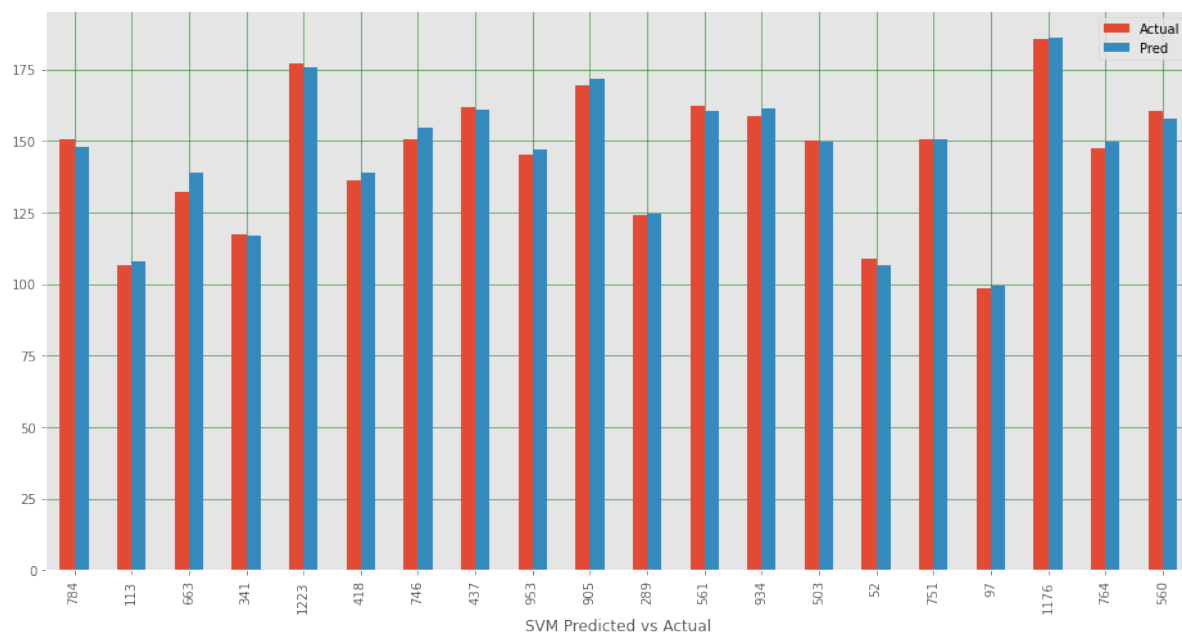
Mean absolute error = 1.44

Mean squared error = 3.86

Median absolute error = 0.97

Explain variance score = 0.99

R2 score = 0.99



In [40]:

```
1 import numpy as np
2 from sklearn.linear_model import LinearRegression
3 from sklearn.svm import SVR
4 from sklearn.model_selection import train_test_split
5
6 # Get the stock data
7 df = df = pd.read_csv('AMZN.csv')
8 # Take a look at the data
9 print(df)
10
11 # Get the Adjusted Close Price
12 df = df[['Adj Close']]
```

	Date	Open	High	Low	Clos
e \					
0	2019-12-02	1804.400024	1805.550049	1762.680054	1781.599976
6					
1	2019-12-03	1760.000000	1772.869995	1747.229980	1769.959961
1					
2	2019-12-04	1774.010010	1789.089966	1760.219971	1760.689941
1					
3	2019-12-05	1763.500000	1763.500000	1740.000000	1740.479980
0					
4	2019-12-06	1751.199951	1754.400024	1740.130005	1751.599976
6					
..
.					
248	2020-11-24	3100.500000	3134.250000	3086.260010	3118.060059
9					
249	2020-11-25	3141.870117	3198.000000	3140.260010	3185.070068
8					
250	2020-11-27	3211.260010	3216.189941	3190.050049	3195.340088
8					
251	2020-11-30	3208.479980	3228.389893	3125.550049	3168.040039
9					
252	2020-12-01	3188.500000	3248.949951	3157.179932	3220.080078
8					
	Adj Close	Volume			
0	1781.599976	3925600			
1	1769.959961	3380900			
2	1760.689941	2670100			
3	1740.479980	2823800			
4	1751.599976	3117400			
..			
248	3118.060059	3602100			
249	3185.070068	3790400			
250	3195.340088	2392900			
251	3168.040039	4063900			
252	3220.080078	4537000			

[253 rows x 7 columns]

In [41]:

```
1 # A variable for predicting 'n' days out into the future
2 forecast_out = 20 #'n=10' days
3 #Create another column (the target ) shifted 'n' units up
4 df['Prediction'] = df[['Adj Close']].shift(-forecast_out)
```

```

4  df['Prediction'] = df[['Adj Close']].shift(-forecast_out)
5  #print the new data set
6  # print(df.tail())
7
8  ### Create the independent data set (X) #####
9  # Convert the dataframe to a numpy array
10 X = np.array(df.drop(['Prediction'],1))
11
12 #Remove the last '30' rows
13 X = X[:-forecast_out]
14 # print(X)
15
16 # ### Create the dependent data set (y) #####
17 # # Convert the dataframe to a numpy array
18 y = np.array(df['Prediction'])
19 # # Get all of the y values except the last '30' rows
20 y = y[:-forecast_out]
21 # print(y)
22
23 # # Split the data into 80% training and 20% testing
24 x_train, x_test, y_train, y_test = train_test_split(X, y, test_
25
26 # # Create and train the Linear Regression Model
27 lr = LinearRegression()
28 # Train the model
29 lr.fit(x_train, y_train)
30
31 # # Create and train the Support Vector Machine (Regressor)
32 svr_rbf = SVR(kernel='rbf', C=1e3, gamma=0.1)
33 svr_rbf.fit(x_train, y_train)
34
35 svm_linear = SVR(kernel='linear',)
36 svm_linear.fit(x_train, y_train)
37
38 svm_confidence2 = svm_linear.score(x_test, y_test)
39 print("linear svm confidence: ", svm_confidence2)
40
41 hnn_amazon=KNeighborsRegressor(n_neighbors = 5)
42 hnn_amazon.fit(x_train,y_train)
43
44 print('Knn Confidence', hnn_amazon.score(x_test, y_test))
45
46 # # Testing Model: Score returns the coefficient of determinati
47 # # The best possible score is 1.0
48 svm_confidence = svr_rbf.score(x_test, y_test)
49 print("rbf svm confidence: ", svm_confidence)
50
51 # # Testing Model: Score returns the coefficient of determinati
52 # # The best possible score is 1.0
53 lr_confidence = lr.score(x_test, y_test)
54 print("lr confidence: ", lr_confidence)

```

linear svm confidence: 0.793673175589624

Knn Confidence 0.9110828900264779

rbf svm confidence: 0.5067684817464231

lr confidence: 0.7902567749251435

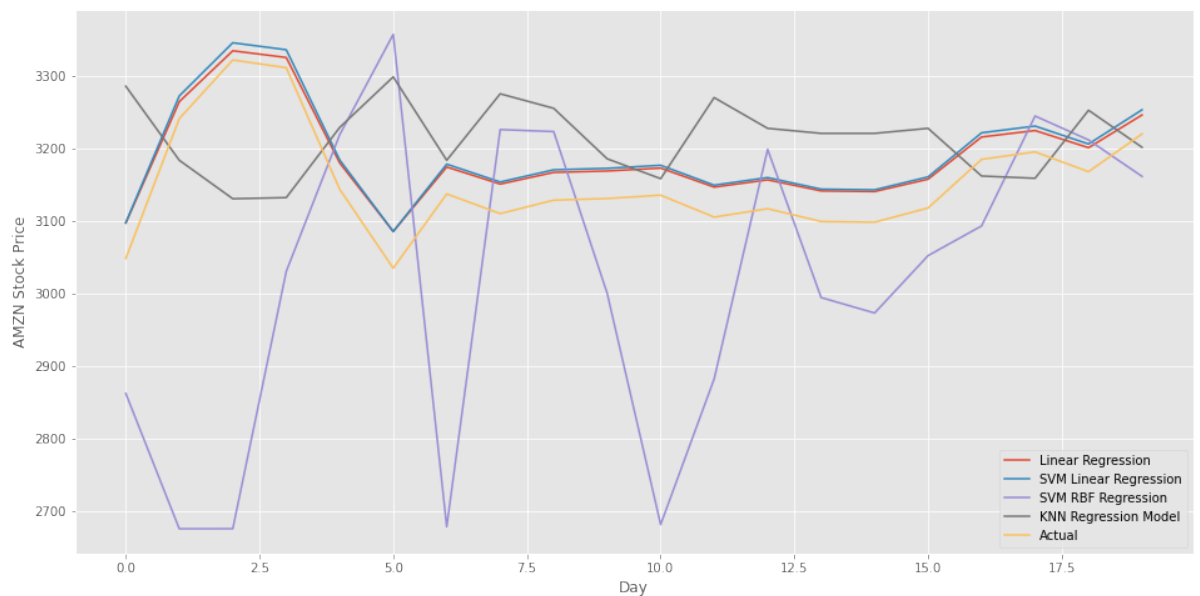
```
In [42]: 1 # Set x_forecast equal to the last 30 rows of the original data
2 x_forecast = np.array(df.drop(['Prediction'],1))[-forecast_out:
3 print(x_forecast)
```

```
[[3048.409912]
 [3241.159912]
 [3322.       ]
 [3311.370117]
 [3143.73999 ]
 [3035.02002 ]
 [3137.389893]
 [3110.280029]
 [3128.810059]
 [3131.060059]
 [3135.659912]
 [3105.459961]
 [3117.02002 ]
 [3099.399902]
 [3098.389893]
 [3118.060059]
 [3185.070068]
 [3195.340088]
 [3168.040039]
 [3220.080078]]
```

```

In [43]: 1 # Print linear regression model predictions for the next '30' d
2 lr_prediction = lr.predict(x_forecast)
3 # print(lr_prediction)
4 # Print support vector regressor model predictions for the next
5 svm_prediction = svr_rbf.predict(x_forecast)
6
7 amazon_knn_pred=hnn_amazon.predict(x_forecast)
8
9 svm_linear_predict = svm_linear.predict(x_forecast)
10
11 actual = df[len(df) - forecast_out:].reset_index()['Adj Close']
12 plot = pd.DataFrame({'Linear Regression': lr_prediction, 'SVM L
13 plot.plot(kind='line', figsize=(16,8))
14 plt.ylabel("AMZN Stock Price")
15 plt.xlabel("Day")
16 plt.show()

```



In []: 1