

IMPORTING LIBRARIES:

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: #here we are loading data using pandas
df=pd.read_csv(r"C:\Users\Hp\Downloads\archive (4) (1).zip")
```

```
In [3]: df
```

Out[3]:	Unnamed: 0	Date	Close	Volume	Open	High	Low
0	0	2024-01-19	2029.3	166078.0	2027.4	2041.9	2022.2
1	1	2024-01-18	2021.6	167013.0	2009.1	2025.6	2007.7
2	2	2024-01-17	2006.5	245194.0	2031.7	2036.1	2004.6
3	3	2024-01-16	2030.2	277995.0	2053.4	2062.8	2027.6
4	4	2024-01-12	2051.6	250946.0	2033.2	2067.3	2033.1
...	...	...	...	...	...	...	...
2506	2528	2014-01-28	1250.5	81426.0	1254.9	1261.9	1248.0
2507	2529	2014-01-27	1263.5	63419.0	1269.9	1280.1	1252.0
2508	2530	2014-01-24	1264.5	34998.0	1264.3	1273.2	1256.9
2509	2531	2014-01-23	1262.5	41697.0	1235.1	1267.1	1230.8
2510	2532	2014-01-22	1238.6	80262.0	1240.5	1243.5	1235.5

2511 rows x 7 columns

Analyzing the each attribute from the above data: Given goldstock prices GoldStockwhich refers to the market where shares of companies involved in the production,exploration,or distribution of gold brought and sold. 1.Date:stock price on particular date. 2.close:At the end of the day,the price of the gold stock per share. 3.Volume:Throughout given date,total shares of the gold stock were bought and sold. 4.Open:when the market opened ,the first trade of the gold stock per share. 5.High:The highest price that the gold stock reached on the particular date. 6.Low:The lowest price that the gold stock reached on the particular date.

```
In [4]: df.head()
```

Out[4]:	Unnamed: 0	Date	Close	Volume	Open	High	Low
0	0	2024-01-19	2029.3	166078.0	2027.4	2041.9	2022.2
1	1	2024-01-18	2021.6	167013.0	2009.1	2025.6	2007.7
2	2	2024-01-17	2006.5	245194.0	2031.7	2036.1	2004.6
3	3	2024-01-16	2030.2	277995.0	2053.4	2062.8	2027.6
4	4	2024-01-12	2051.6	250946.0	2033.2	2067.3	2033.1

```
In [5]: df.tail()
```

Out[5]:	Unnamed: 0	Date	Close	Volume	Open	High	Low
2506	2528	2014-01-28	1250.5	81426.0	1254.9	1261.9	1248.0
2507	2529	2014-01-27	1263.5	63419.0	1269.9	1280.1	1252.0
2508	2530	2014-01-24	1264.5	34998.0	1264.3	1273.2	1256.9
2509	2531	2014-01-23	1262.5	41697.0	1235.1	1267.1	1230.8
2510	2532	2014-01-22	1238.6	80262.0	1240.5	1243.5	1235.5

```
In [6]: df.isnull().sum()
```

Out[6]:	Unnamed: 0	Date	Close	Volume	Open	High	Low
Unnamed: 0	0						
Date	0						
Close	0						
Volume	0						
Open	0						
High	0						
Low	0						
dtype:	int64						

```
In [7]: df.duplicated().sum()
```

```
Out[7]: 0
```

```
In [8]: df['Date']=pd.to_datetime(df['Date'])
df['dayofweek']=df['Date'].dt.dayofweek
df['Month']=df['Date'].dt.month
df['Year']=df['Date'].dt.year
```

```
In [9]: df1 = df.drop(df.columns[1], axis=1)
df1
```

Out[9]:	Unnamed: 0	Close	Volume	Open	High	Low	dayofweek	Month	Year
0	0	2029.3	166078.0	2027.4	2041.9	2022.2	4	1	2024
1	1	2021.6	167013.0	2009.1	2025.6	2007.7	3	1	2024
2	2	2006.5	245194.0	2031.7	2036.1	2004.6	2	1	2024
3	3	2030.2	277995.0	2053.4	2062.8	2027.6	1	1	2024
4	4	2051.6	250946.0	2033.2	2067.3	2033.1	4	1	2024
...	...	...	...	...	...	...	...	...	...
2506	2528	1250.5	81426.0	1254.9	1261.9	1248.0	1	1	2014
2507	2529	1263.5	63419.0	1269.9	1280.1	1252.0	0	1	2014
2508	2530	1264.5	34998.0	1264.3	1273.2	1256.9	4	1	2014
2509	2531	1262.5	41697.0	1235.1	1267.1	1230.8	3	1	2014
2510	2532	1238.6	80262.0	1240.5	1243.5	1235.5	2	1	2014

2511 rows x 9 columns

```
In [10]: df2 = df1.drop(df1.columns[0], axis=1)
df2
```

Out[10]:	Close	Volume	Open	High	Low	dayofweek	Month	Year
0	2029.3	166078.0	2027.4	2041.9	2022.2	4	1	2024
1	2021.6	167013.0	2009.1	2025.6	2007.7	3	1	2024
2	2006.5	245194.0	2031.7	2036.1	2004.6	2	1	2024
3	2030.2	277995.0	2053.4	2062.8	2027.6	1	1	2024
4	2051.6	250946.0	2033.2	2067.3	2033.1	4	1	2024
...	...	...	...	...	...	...	...	...
2506	1250.5	81426.0	1254.9	1261.9	1248.0	1	1	2014
2507	1263.5	63419.0	1269.9	1280.1	1252.0	0	1	2014
2508	1264.5	34998.0	1264.3	1273.2	1256.9	4	1	2014
2509	1262.5	41697.0	1235.1	1267.1	1230.8	3	1	2014
2510	1238.6	80262.0	1240.5	1243.5	1235.5	2	1	2014

2511 rows x 8 columns

```
In [11]: Q1 = df2.quantile(0.25)
Q3 = df2.quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
data_cleaned = df2[~((df2 < lower_bound) | (df2 > upper_bound)).any(axis=1]]

print("Original Data:")
print(df2)

print("\nCleaned Data:")
print(data_cleaned)
```

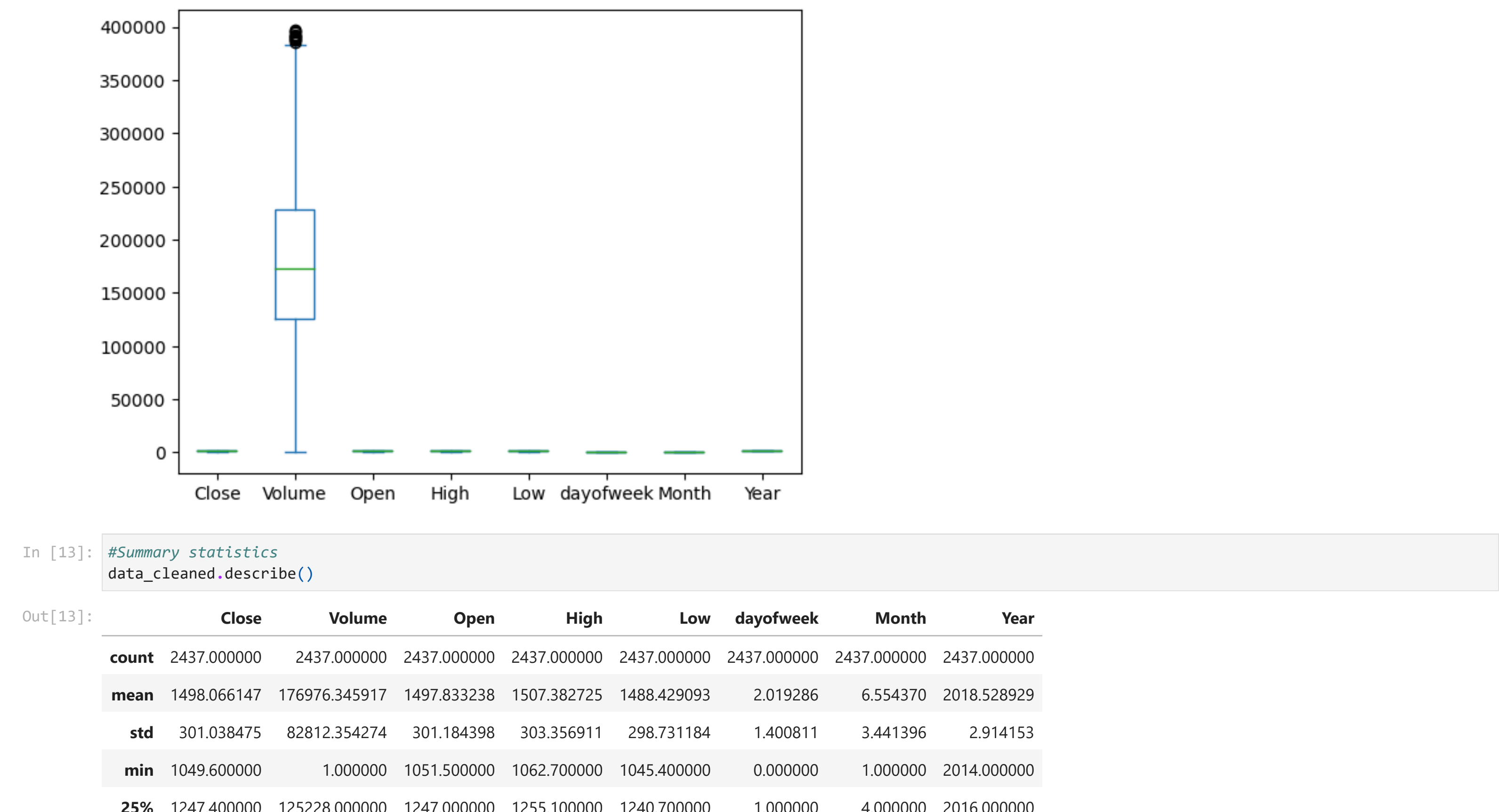
Original Data:	Close	Volume	Open	High	Low	dayofweek	Month	Year
0	2029.3	166078.0	2027.4	2041.9	2022.2	4	1	2024
1	2021.6	167013.0	2009.1	2025.6	2007.7	3	1	2024
2	2006.5	245194.0	2031.7	2036.1	2004.6	2	1	2024
3	2030.2	277995.0	2053.4	2062.8	2027.6	1	1	2024
4	2051.6	250946.0	2033.2	2067.3	2033.1	4	1	2024
...	...	...	...	...	...	...	...	...
2506	1250.5	81426.0	1254.9	1261.9	1248.0	1	1	2014
2507	1263.5	63419.0	1269.9	1280.1	1252.0	0	1	2014
2508	1264.5	34998.0	1264.3	1273.2	1256.9	4	1	2014
2509	1262.5	41697.0	1235.1	1267.1	1230.8	3	1	2014
2510	1238.6	80262.0	1240.5	1243.5	1235.5	2	1	2014

[2511 rows x 8 columns]

Cleaned Data:	Close	Volume	Open	High	Low	dayofweek	Month	Year
0	2029.3	166078.0	2027.4	2041.9	2022.2	4	1	2024
1	2021.6	167013.0	2009.1	2025.6	2007.7	3	1	2024
2	2006.5	245194.0	2031.7	2036.1	2004.6	2	1	2024
3	2030.2	277995.0	2053.4	2062.8	2027.6	1	1	2024
4	2051.6	250946.0	2033.2	2067.3	2033.1	4	1	2024
...	...	...	...	...	...	...	...	...
2506	1250.5	81426.0	1254.9	1261.9	1248.0	1	1	2014
2507	1263.5	63419.0	1269.9	1280.1	1252.0	0	1	2014
2508	1264.5	34998.0	1264.3	1273.2	1256.9	4	1	2014
2509	1262.5	41697.0	1235.1	1267.1	1230.8	3	1	2014
2510	1238.6	80262.0	1240.5	1243.5	1235.5	2	1	2014

[2437 rows x 8 columns]

```
In [12]: data_cleaned.plot(kind='box')
```



```
In [13]: #Summary statistics
data_cleaned.describe()
```

Out[13]:	Close	Volume	Open	High	Low	dayofweek	Month	Year
count	2437.000000	2437.000000	2437.000000	2437.000000	2437.000000	2437.000000	2437.000000	2437.000000
mean	1498.066147	176976.345917	1497.833238	1507.382725	1488.429093	2.019286	6.554370	2018.528929
std	301.038475	82812.354274	301.184398	303.356911	298.731184	1.400811	3.441396	2.914153
min	1049.600000	1.000000	1051.500000	1062.700000	1045.400000	0.000000	1.000000	2014.000000
25%	1247.400000	125228.000000	1247.000000	1255.100000	1240.700000	1.000000	4.000000	2016.000000
50%	1329.500000	172801.000000	1330.800000	1337.600000	1323.600000	2.000000	7.000000	2018.000000
75%	1807.200000	228720.000000	1807.400000	1816.500000	1794.100000	3.000000	10.000000	2021.000000
max	2093.100000	396657.000000	2094.400000	2098.200000	2074.600000	4.000000	12.000000	2024.000000

```
In [14]: #correlation matrix
correlation_matrix=data_cleaned.corr()
correlation_matrix
```

Out[14]:	Close	Volume	Open	High	Low	dayofweek	Month	Year
Close	1.000000	0.019988	0.999130	0.999623	0.999616	0.008208	0.013139	0.893486
Volume	0.019988	1.000000	0.022064	0.024467	0.017113	0.050983	-0.005889	0.104433
Open	0.999130	0.022064	1.000000	0.999536	0.999523	0.006848	0.013119	0.893521
High	0.999623	0.024467	0.999536	1.000000	0.999426	0.008177	0.012042	0.892439
Low	0.999616	0.017113	0.999523	0.999426	1.000000	0.007288	0.014914	0.894628
dayofweek	0.008208	0.050983	0.006848	0.008177	0.007288	1.000000	-0.017887	0.008763
Month	0.013139	-0.005889	0.013119	0.012042	0.014914	-0.017887	1.000000	-0.023847
Year	0.893486	0.104433	0.893521	0.892439	0.894628	0.008763	-0.023847	1.000000

## Building an algorithm

```
In [15]: #importing some libraries to build a model
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error as mse
from sklearn.metrics import r2_score
```

```
In [16]: x=df[['Volume','Open','High','Low','dayofweek','Month','Year']]
y=df[['Close']]
x_train,x_test,y_train,y_test= train_test_split(x,y,test_size=0.2,random_state=49)
```

Here we choose target variable as 'close'. we choose 'Volume','Open','High','Low','dayofweek','Month','Year'variables as input variables.

```
In [17]: x
```

Out[17]:	Volume	Open	High	Low	dayofweek	Month	Year
0	166078.0	2027.4	2041.9	2022.2	4	1	2024
1	167013.0	2009.1	2025.6	2007.7	3	1	2024
2	245194.0	2031.7	2036.1	2004.6	2	1	2024
3	277995.0	2053.4	2062.8	2027.6	1	1	2024
4	250946.0	2033.2	2067.3	2033.1	4	1	2024
...	...	...	...	...	...	...	...
2506	81426.0	1254.9	1261.9	1248.0	1	1	2014
2507	63419.0	1269.9	1280.1	1252.0	0	1	2014
2508	34998.0	1264.3	1273.2	1256.9	4	1	2014
2509	41697.0	1235.1	1267.1	1230.8	3	1	2014
2510	80262.0	1240.5	1243.5	1235.5	2	1	2014

2511 rows x 7 columns

```
In [18]: y
```

Out[18]:	Close
0	2029.3
1	2021.6
2	2006.5
3	2030.2
4	2051.6
...	...
2506	1250.5
2507	1263.5
2508	1264.5
2509	1262.5
2510	1238.6

2511 rows x 1 columns

```
In [19]: from sklearn.preprocessing import MinMaxScaler
from sklearn.preprocessing import StandardScaler
#feature Scaling
scaler =StandardScaler().fit(x_train)
```

```
In [20]: scaler
```

```
Out[20]: StandardScaler
StandardScaler()
```

```
In [21]: #creating a linear model
from sklearn.linear_model import LinearRegression
lm=LinearRegression()
lm.fit(x_train,y_train)
```

```
Out[21]: LinearRegression
LinearRegression()
```

```
In [22]: y_pred=lm.predict(x_test)
len(y_train)
```

```
Out[22]: 2808
```

```
In [23]: # predictions on the training data
y_train_pred = lm.predict(x_train)
```

```
# Plotting the actual vs predicted values for the training data
plt.figure(figsize
```