

The Price of Gold

1 Domain-specific area and objectives of the project

The chosen domain for this report is financial markets, in particular, tracking one of the most popular commodity and a form of currency, **gold**.

Gold has been used as a unit of currency and a store of value for centuries[1]. Its historical significance as a medium of exchange as well as symbol of wealth has contributed to its enduring role in finance. Gold is actively traded on daily basis on various global exchanges. Traders buy and sell gold futures contracts to speculate on its future price movements.

Gold is incorporated into various financial instruments, such as Exchange-Traded Funds (ETFs) and gold-backed securities[2]. These instruments allow commercial and retail investors to gain exposure to gold prices without physically owning the metal. It also helps to fight against inflation, as gold supply cannot be artificially inflated like fiat currencies[3].

As well as being a financial instrument to invest in and generate wealth, gold is often considered a "safe haven" asset. During times of economic uncertainty or political instability, people tend to flee towards gold as a store of value. Why? Gold is perceived as a stable and reliable asset that can retain its worth in challenging economic conditions. It can be said that people's perception and ideologies towards gold contributes to the metal's value as well.

As of now, with global market instability, war, and inflation determining where the price of gold will be in the near future can be a problem for both retail investor and institutions alike.

Now that we know what gold is and why people value it, we can understand why Machine Learning (ML) can come into place here.

A linear regression model can help us predict the future price of gold. It can give investors an indicator of whether their asset will appreciate or depreciate in the coming future. It can also provide a general consensus of the public's view on gold, as declining gold prices could indicate a disinterest in gold as a commodity, as its demand may be dropping. *Though there are several factors the effect the price of gold, using historical gold price data is a valid method of predicting its future price.* There are already algorithms and trading bots being used by traders that utilize ML to make trades based of historical data.

The objective of this project is to build a model that will give us a prediction to future gold prices, which can be used to aid in our decision-making process of whether to add gold to our investment portfolio, or to liquidate it, as well as visualizing the price change over the decades.

The results of this project may offer both academic and practical insights into the precious metals markets.

2 Dataset

The dataset required to be used for this project comes from Kaggle[4]. The dataset comes in three CSV files that are zipped together. The file names are "Gold_Daily", "Gold_Monthly", and "Gold_Yearly". The size of each files are 358Kb, 23Kb and 3Kb, bringing the total size of the dataset to 384Kb. The way the data was sourced, as stated on the Kaggle page, was from Investing.com[5], which is financial market website that also stores historical prices of tradable assets on the exchange.

The three CSV files originally contains information such as the date of the historical gold price, its high, low, open, as well as its volume and change. The total number of columns for each file is 7. The date column is the "date" of the historical gold price, the "price" is its closing price for that day. The "open" column is what the price of gold opened at during the start of the day, and the "high" and "low" is the highest price and lowest price gold reached during the day. The original dataset contains missing values and doesn't have a unique ID for each row.

In the mid-term webinar hosted by Dr Georgios Matorakis, he make mentions it is allowed to add random noise, remove random values and intentionally make the data set dirty, in order to make full use of the pre-processing step. In addition to the original dataset already containing missing values, I will be adding random noise into the values, like multiple values in cells.

This dataset is suitable to the objectives of the project as it contains a linear trendline that can be implemented into a linear regression model. It highlights the chosen domain of financial markets, and the tracking price of one of the most popular commodities.

3 Data Preparation

We will start processing the dataset here. The steps that will be undertaken here are:

- converting into appropriate datatype
- addressing missing values
- assigning a unique ID to the dataset
- transforming it into a 1NF dataset

These steps are needed to reduce redundancy, increase accuracy of the dataset, which will in turn help our linear regression model later on. We will be using the Pandas library for this.

```
In [1]: import pandas as pd

#load the 3 csv files into a dataframe
daily_df = pd.read_csv('gold/Gold_Daily.csv')
monthly_df = pd.read_csv('gold/Gold_Monthly.csv')
yearly_df = pd.read_csv('gold/Gold_Yearly.csv')

print(daily_df.head())
print("\n")
print(monthly_df.head(10))
print("\n")
print(yearly_df.head())
print("\n")

      Date  Price  Open  High  Low  Vol  Change %
0  20-Dec-94  375.8  375.6  376.5  374.8  0.008  -0.008
1  07-Dec-94  376.1  377.0  377.5  375.1  1.89K  0.008
2  08-Dec-94  376.6  375.3  376.7  374.8  0.48K  0.019
3  09-Dec-94  377.0  375.6  377.0  375.5  0.08K  0.011
4  12-Dec-94  377.5  376.8  377.8  376.5  0.37K  0.013

      Date  Price  Open  High  Low  Vol  Change %
0  Feb-79  271.6  271.6  271.6  271.6  0.0  0.0727
1  Mar-79  256.6  256.6  256.6  256.6  0.0  -0.0552
2  Apr-79  262.4  262.4  262.4  262.4  0.0  0.0226
3  May-79  290.7  290.7  290.7  290.7  0.0  0.1079
4  Jun-79  293.5  293.5  293.5  293.5  0.0  0.0096
5  Jul-79  297.3  297.3  297.3  297.3  297.3  0.0129
6  Aug-79  326.7  326.7  326.7  326.7  0.0  0.0989
7  Sep-79  404.5  404.5  404.5  404.5  0.0  0.2381
8  Oct-79  385.7  385.7  385.7  385.7  0.0  -0.0465
9  Nov-79  419.1  396.3  421  390.0  0.0  0.0866

Year  YearAvg  YearClosing  Price  Year  Year  High  Year  Low  Year  Close  \
0  1969          41.10         41.80         43.75         35.00         35.21
1  1970          35.96         35.13         39.19         34.78         37.38
2  1971          40.80         40.80         43.90         37.33         43.50
3  1972          58.17         43.73         70.00         43.73         64.70
4  1973          97.12         64.99         127.00         64.10         112.25

Annual % Change
0          -0.1607
1           0.0616
2           0.1637
3           0.4874
4           0.7349

In [2]: daily_df.info()
print("\n")
monthly_df.info()
print("\n")
yearly_df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6886 entries, 0 to 6885
Data columns (total 7 columns):
 #   Column      Non-Null Count  Dtype  ---
 0   Date       6886 non-null     object
 1   Price      6886 non-null     float64
 2   Open       6886 non-null     float64
 3   High       6886 non-null     float64
 4   Low        6886 non-null     float64
 5   Vol        6886 non-null     object
 6   Change %   6886 non-null     float64
dtypes: float64(6), object(1)
memory usage: 376.7+ KB

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 515 entries, 0 to 514
Data columns (total 7 columns):
 #   Column      Non-Null Count  Dtype  ---
 0   Date       515 non-null     object
 1   Price      515 non-null     float64
 2   Open       515 non-null     float64
 3   High       515 non-null     object
 4   Low        515 non-null     float64
 5   Vol        515 non-null     object
 6   Change %   515 non-null     float64
dtypes: float64(6), object(1)
memory usage: 28.3+ KB

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 53 entries, 0 to 52
Data columns (total 7 columns):
 #   Column      Non-Null Count  Dtype  ---
 0   Date       53 non-null     object
 1   Price      53 non-null     float64
 2   Open       53 non-null     float64
 3   High       53 non-null     float64
 4   Low        53 non-null     float64
 5   Vol        53 non-null     float64
 6   Annual %   53 non-null     float64
dtypes: float64(6), int64(1)
memory usage: 3.0 KB
```

The above shows the format of the three CSV files, as well as the information of each dataframe. We can see they all contain 7 columns, with both float and object(string/mixed) data types. Lets first address the missing values for each file.

In the daily gold price csv file, the volume column has missing volume data denoted with the symbol '.': Lets replace it with nan first To do that, we will need to import numpy

```
In [3]: daily_df[10:20]

Out[3]:
      Date  Price  Open  High  Low  Vol  Change %
10 20-Dec-94  381.8  381.8  382.3  381.4  0.02K  0.009
11 21-Dec-94  381.6  382.0  382.0  381.5  0.03K  -0.005
12 22-Dec-94  381.0  381.5  381.5  381.5  0.05K  -0.016
13 23-Dec-94  380.4  381.0  381.0  380.4  0.02K  -0.006
14 27-Dec-94  381.6  381.0  382.0  381.0  0.06K  0.003
15 28-Dec-94  383.7  382.1  383.6  382.1  - 0.0052
16 29-Dec-94  382.7  382.7  382.7  382.7  - -0.0026
17 30-Dec-94  383.1  383.1  383.1  383.1  - 0.0001
18 03-Jan-95  379.6  379.6  379.6  379.6  - -0.0091
19 04-Jan-95  374.0  377.7  377.7  374.0  0.00K  -0.0148

In [4]: #identify duplicate rows
duplicateRows = monthly_df[daily_df.duplicated()]

#view duplicate rows
duplicateRows_D

Out[4]:
      Date  Price  Open  High  Low  Vol  Change %

No duplicate rows

In [5]: import numpy as np
#daily gold price
#replace '-' with NaN
daily_df['Vol'] = daily_df['Vol'].replace('-', np.nan)
print(daily_df[10:20])

#count number of nan values
daily_df['Vol'].isna().sum()

      Date  Price  Open  High  Low  Vol  Change %
10 20-Dec-94  381.8  381.8  382.3  381.4  0.02K  0.009
11 21-Dec-94  381.6  382.0  382.0  381.5  0.03K  -0.005
12 22-Dec-94  381.0  381.5  381.5  381.5  0.05K  -0.0016
13 23-Dec-94  380.4  381.0  381.0  380.4  0.02K  -0.006
14 27-Dec-94  381.6  381.0  382.0  381.0  0.06K  0.003
15 28-Dec-94  383.7  382.1  383.6  382.1  NaN  0.0055
16 29-Dec-94  382.7  382.7  382.7  382.7  NaN  -0.0026
17 30-Dec-94  383.1  383.1  383.1  383.1  NaN  0.0010
18 03-Jan-95  379.6  379.6  379.6  379.6  NaN  -0.0091
19 04-Jan-95  374.0  377.7  377.7  374.0  0.00K  -0.0148
1302
```

We have replaced the missing values for the volume column with numpy NaN, and we can see that we have a total missing value count of 1302. We need to address this in our process of normalizing this dataset as it will improve data integrity and model accuracy.

How to address missing values?

Because our dataset is linear, we can use interpolate() to fill the missing values. The volume column denoted the volume traded in the session in thousands. We need to remove the "K" from the object and convert it to a float first, before interpolation can begin.

For columns that are identical with the daily and monthly csv files, we will be performing functions so that we can call upon them to clean the data set for the "Gold_Monthly.csv" file

```
In [6]: #function to remove 'K' string and convert to numeric
def convertVol(dataframe):
    #remove 'K' and convert to numeric
    dataframe['Vol'] = dataframe['Vol'].str.replace('K', '').astype(float)
    #multiply by 1000 to convert from 'str' (thousands) to actual numeric values
    dataframe['Vol'] = dataframe['Vol'] * 1000
    return dataframe['Vol']

daily_df['Vol'] = convertVol(daily_df)

print(daily_df[10:20])

      Date  Price  Open  High  Low  Vol  Change %
10 20-Dec-94  381.8  381.8  382.3  381.4  0.02K  0.009
11 21-Dec-94  381.6  382.0  382.0  381.5  0.03K  -0.005
12 22-Dec-94  381.0  381.5  381.5  381.5  0.05K  -0.0016
13 23-Dec-94  380.4  381.0  381.0  380.4  0.02K  -0.006
14 27-Dec-94  381.6  381.0  382.0  381.0  0.06K  0.003
15 28-Dec-94  383.7  382.1  383.6  382.1  NaN  0.0055
16 29-Dec-94  382.7  382.7  382.7  382.7  NaN  -0.0026
17 30-Dec-94  383.1  383.1  383.1  383.1  NaN  0.0010
18 03-Jan-95  379.6  379.6  379.6  379.6  NaN  -0.0091
19 04-Jan-95  374.0  377.7  377.7  374.0  0.00K  -0.0148

In [7]: #interpolate the missing values, set the direction to forward, from first entry to last
daily_df = daily_df.interpolate(method='linear', limit_direction='forward')

#convert date to pandas datetime format
daily_df['Date'] = pd.to_datetime(daily_df['Date']).format('%d-%b-%y')

#rename change column
daily_df = daily_df.rename(columns={"Change %": "Change"})

#save the date column values to a variable, so we can plot it later
daily_Dates_List = daily_df['Date'].tolist()

#assign unique ID for each row
daily_df['ID'] = range(len(daily_df))

#rearrange columns
arrange = ["ID", "Date", "Price", "Open", "High", "Low", "Vol", "Change"]
daily_df = daily_df.reindex(columns=arrange)

print(daily_df[10:10])
print("\n")
print("Count of missing values:", daily_df['Vol'].isna().sum())

      Date  Price  Open  High  Low  Vol  Change %
0  1994-12-06  376.8  375.6  376.5  375.1  990.0  0.008
1  1994-12-07  376.1  377.0  377.5  375.1  1880.0  0.008
2  1994-12-08  376.6  375.3  376.7  374.8  480.0  0.019
3  1994-12-09  377.0  375.6  377.0  375.5  80.0  0.011
4  1994-12-10  377.5  376.8  377.8  376.5  70.0  0.013
5  1994-12-13  377.6  379.5  379.5  379.0  40.0  0.026
6  1994-12-14  378.9  379.0  381.0  379.0  50.0  0.034
7  1994-12-15  380.2  379.9  380.3  378.9  80.0  0.034
8  1994-12-16  378.6  380.5  380.5  379.0  40.0  0.042
9  1994-12-19  379.2  379.2  379.2  379.2  10.0  0.016

Count of missing values: 0
```

```
In [8]: daily_df.dtypes

ID          int32
Date       datetime64[ns]
Price      float64
Open       float64
High       float64
Low        float64
Vol        float64
Change     float64
dtype: object
```

We can see the volume column has been converted into a float, and the missing values has been interpolated. The total missing values is now zero.

I have also converted the "Date" column to the appropriate pandas datetime object, as well as renaming the "Change%" column, which is the percent change of the previous and current (day, month, or year) price, to "Change", to remove the special characters.

Lets now proceed with prepping the "Gold_Monthly.csv". As mentioned above, this dataset has noise introduced into it with some cells having values.

```
In [9]: #identify duplicate rows
duplicateRows = monthly_df[monthly_df.duplicated()]

#view duplicate rows
duplicateRows

Out[9]:
      Date  Price  Open  High  Low  Vol  Change %

In [10]: #rename columns
monthly_df = monthly_df.rename(columns={"Vol": "Vol", "Change %": "Change"})

#replace '-' with NaN
monthly_df['Vol'] = monthly_df['Vol'].replace('-', np.nan)

#the 2nd last entry in monthly where the volume is in millions, we will convert that first
monthly_df['Vol'] = monthly_df['Vol'].str.replace('M', '')
monthly_df['Vol'] = monthly_df['Vol'] * 1.31 * 1000000

#remove 'K' and convert to float
monthly_df['Vol'] = convertVol(monthly_df)

In [11]: #fill missing values
monthly_df = monthly_df.interpolate(method='linear', limit_direction='forward')

#get pandas date time
monthly_df['Date'] = pd.to_datetime(monthly_df['Date']).format('%b-%y')

#save the date column values to a variable, so we can plot it later
daily_Dates_List = monthly_df['Date'].tolist()

#remove the first 11 dates, cuz we remove the NaN values later
for x, item in enumerate(monthly_Dates_List):
    if x == 10:
        break

print(monthly_df[0:20])
print("\n")
print("Count of missing values:", monthly_df['Vol'].isna().sum())

      Date  Price  Open  High  Low  Vol  Change %
0  1979-02-01  271.6  271.6  271.6  271.6  NaN  0.0727
1  1979-03-01  256.6  256.6  256.6  256.6  NaN  -0.0552
2  1979-04-01  262.4  262.4  262.4  262.4  NaN  0.0226
3  1979-05-01  290.7  290.7  290.7  290.7  NaN  0.1079
4  1979-06-01  293.5  293.5  293.5  293.5  NaN  0.0096
5  1979-07-01  297.3  297.3  297.3  297.3  NaN  0.0129
6  1979-08-01  326.7  326.7  326.7  326.7  NaN  0.0989
7  1979-09-01  404.5  404.5  404.5  404.5  NaN  0.2381
8  1979-10-01  385.7  385.7  385.7  385.7  NaN  -0.0465
9  1979-11-01  419.1  396.3  421  390.0  NaN  0.0866
10 1979-12-01  533.6  430.5  534.5  424.0  NaN  0.2732
11 1980-01-01  681.5  562.5  875.875  558.0  17930.0  0.272
12 1980-02-01  631.0  677.0  729  599.0  4930.0  -0.0741
13 1980-03-01  501.5  631.0  648.5  453.0  5700.0  -0.2052
14 1980-04-01  582.0  504.0  562  465.0  2330.0  0.0002
15 1980-05-01  545.2  490.6  552  478.0  2570.0  0.0869
16 1980-06-01  647.4  567.0  658.5  550.0  2651.0  0.1875
17 1980-07-01  619.7  662.9  691  604.0  3130.0  -0.0428
18 1980-08-01  635.0  620.0  649  600.0  14270.0  0.0247
19 1980-09-01  671.5  639.7  727  635.5  48930.0  0.0575

Count of missing values: 11
```

We have 3 cells in our dataset that has commas, which indicate multiple values. Lets remove the value of the before the comma, as well as the comma function.

```
In [13]: #lambda function to each cell in the specified column, split the string, and keep the second element only
monthly_df['High'] = monthly_df['High'].apply(lambda x: x.split(',')[1] if ',' in x else x)
print(monthly_df[10:20])

      Date  Price  Open  High  Low  Vol  Change %
0  271.6
1  256.6
2  262.4
3  290.7
4  293.5
5  297.3
6  326.7
7  404.5
8  385.7
9  421
10 534.5
11 875
12 729
13 648.5
14 562
15 552
16 658.5
17 691
18 649
19 727
Name: High, dtype: object
```

```
In [14]: #convert dtype object to float
monthly_df = monthly_df.astype({'High': 'float64'})
print(monthly_df.dtypes)

Date          datetime64[ns]
Price         float64
Open          float64
High          float64
Low           float64
Vol           float64
Change        float64
dtype: object
```

Looking at the above, we can see our data processing for the "Monthly_Gold.csv" file gave us 11 NaN values. This is because, originally, the first 11 rows had no recordings of the monthly volume of transactions, and our interpolate functions that fills the NaN value starts from the first index. Since there is no recorded, the interpolation can only being from the first recording of volume, in this case, row ID 11.

You might say, why not interpolate backwards from the last index? This would not make sense for this data set as logically speaking, the volume of transactions recorded at the year 2021, would naturally be significantly higher than that of 1979. Interpolating the data backwards would give us highly inaccurate estimates of the missing values.

The logical step here would be to drop the 11 rows.

```
In [15]: monthly_df = monthly_df.dropna()

#assign unique ID for each row
monthly_df['ID'] = range(len(monthly_df))

#rearrange columns
arrange = ["ID", "Date", "Price", "Open", "High", "Low", "Vol", "Change"]
monthly_df = monthly_df.reindex(columns=arrange)

print(monthly_df[10:10])

Out[15]:
      ID      Date  Price  Open  High  Low  Vol  Change %
11 0 1980-01-01  681.5  562.5  875.0  558.0  17990.0  0.272
12 1 1980-02-01  631.0  670.0  729.0  599.0  4930.0  -0.0741
13 2 1980-03-01  501.5  630.0  648.5  453.0  5700.0  -0.2052
14 3 1980-04-01  581.5  504.0  562.0  465.0  2390.0  0.0002
15 4 1980-05-01  545.2  490.6  552.0  478.0  2570.0  0.0869
```

Now prepping the "Yearly_Gold.csv". Only changes here needed are assigning a unique ID, and renaming of columns.

```
In [16]: #rename columns
yearly_df = yearly_df.rename(columns={"Average": "Average", "Annual % Change": "Change"})

#assign ID
yearly_df['ID'] = range(len(yearly_df))

#arrange columns
cols = list(yearly_df.columns)
cols = [cols[-1]] + cols[:-1]
yearly_df = yearly_df[cols]

print(yearly_df.head())

Out[16]:
      ID  Year  Average  Year Open  Year High  Year Low  Year Close  Change
0  0  1969          41.10         41.80         43.75         35.00         35.21  -0.1607
1  1  1970          35.96         35.13         39.19         34.78         37.38   0.0616
2  2  1971          40.80         37.33         43.90         37.33         43.50   0.1637
3  3  1972          58.17         43.73         70.00         43.73         64.70   0.4874
4  4  1973          97.12         64.99         127.00         64.10         112.25   0.7349
```

All three dataset have been normalized and converted in 1NF, where there following are satisfied:

- a single cell must not hold more than one value (atomicity)
- there must be a primary key for identification
- no duplicated rows or columns
- each column must have only one value for each row in the table

4 Statistical Analysis

Our cleaned dataset can now be analyzed to provide statistical insights. These statistical insights removes unnecessary information and logs important data about our dataset in a succinct manner.

These will include:

- measures of central tendency
- measures of spread
- type of distribution

4.1 Measures of central tendency

The mode, mean and median of the dataset.

```
In [17]: import matplotlib.pyplot as plt
import seaborn as sns

#display mode mean and median of monthly and daily gold prices
def measureCentralTendency(dataframe, column):
    if column == 'Price':
        measure_list = []
        measure_list.append(dataframe['Price'].mode()[0])
        measure_list.append(dataframe['Price'].mean())
        measure_list.append(dataframe['Price'].median())
        return measure_list
    if column == 'Average':
        measure_list = []
        measure_list.append(dataframe['Average'].mode()[0])
        measure_list.append(dataframe['Average'].mean())
        measure_list.append(dataframe['Average'].median())
        measure_list.append(dataframe['Average'].std())
        return measure_list

getStatList = [measureCentralTendency(daily_df, 'Price'), measureCentralTendency(monthly_df, 'Price'),
               measureCentralTendency(yearly_df, 'Average')]

stats = pd.DataFrame(getStatList, columns=["Mode", "Mean", "Median", "Standard Deviation"],
                     index=["Daily Price", "Monthly Price", "Yearly Price"])
stats

Gold Price Statistics 1979-2021

Out[18]:
      Mode      Mean      Median  Standard Deviation
Daily Price  382.30  909.824807  890.90      541.516040
Monthly Price  342.70  723.299107  419.20      498.075312
Yearly Price   35.96  591.265094  383.73      497.289713
```

If we look at the daily average price of gold, we can see its around 909, with a standard deviation of 541. Why is the standard deviation so large? The standard deviation is the measure of spread against the mean. For gold prices, it makes sense as to why the sample values in the dataset, meaning its value is very so much from the mean. Gold is an appreciating and holding asset that has been rising over the decades, meaning its price has been going up for the most part, linearly. So if the average daily price of gold has been around 909, having lower gold price value in the early 90s, and higher values in the 2000s when the asset appreciated, explains the high standard deviation against the mean. The mean is heavily influenced by large outliers.

4.2 Measures of Spread

Standard Deviation and Skewness

```
In [19]: mean_price = daily_df['Price'].mean()
std_dev_price = daily_df['Price'].std()

#the mean as a horizontal line
plt.axhline(mean_price, color='red', linestyle='dashed', linewidth=2, label='Mean Price')

#the standard deviation as error bars
plt.errorbar(0, mean_price, std_dev_price, fmt='o', color='blue', capsize=5, label='Standard Deviation')

#add labels and title
plt.title('Daily Gold')
plt.xlabel('Daily Price')
plt.legend()

plt.show()

Daily Gold

Out[20]:
fig, axes = plt.subplots(1, 3, figsize=(15, 3))

axes[0].title.set_text('Daily')
axes[0].boxplot(daily_df['Price'], vert=False)

axes[1].title.set_text('Monthly')
axes[1].boxplot(monthly_df['Price'], vert=False)

axes[2].title.set_text('Yearly')
axes[2].boxplot(yearly_df['Year Open'], vert=False)

plt.show()

Daily Monthly Yearly
```

Some outliers in yearly opening price. We do not remove them as in this case, they are still relevant to the overall objective.

right skewed/positively skewed distribution of volume

```
In [21]: #daily and monthly volume histograms
fig, axes = plt.subplots(1, 2, figsize=(12, 3))

axes[0].hist(daily_df['Vol'])
axes[1].hist(monthly_df['Vol'])

plt.show()

Daily Volume States
count 6886.000000
mean 2202.091926
std 10502.164473
min 0.000000
25% 20.000000
50% 80.000000
75% 357.000000
max 242860.000000
Name: Vol, dtype: float64

Monthly Volume States
count 515.0000
mean 20688.43254
std 24645.10263
min 0.000000
25% 4150.00000
50% 10990.00000
75% 19900.00000
max 179990.00000
Name: Vol, dtype: float64
```

The dataset on daily gold transaction volumes comprises 6,886 data points. The mean, or average, daily volume of gold transactions is calculated at approximately 2202, with a standard deviation of 10502, indicating a notable degree of variability in the data.

The recorded volumes range from a minimum of 0 to a maximum of 242,860, signifying a diverse distribution. Quartile analysis reveals that 25% of the daily volumes fall below 20 (Q1), the first quartile, 50% fall below 80 (the median), and 75% fall below 357.5 (Q3, the third quartile).

The large difference between the 75th percentile and the maximum value suggests the presence of outliers or days with exceptionally high transaction volumes.

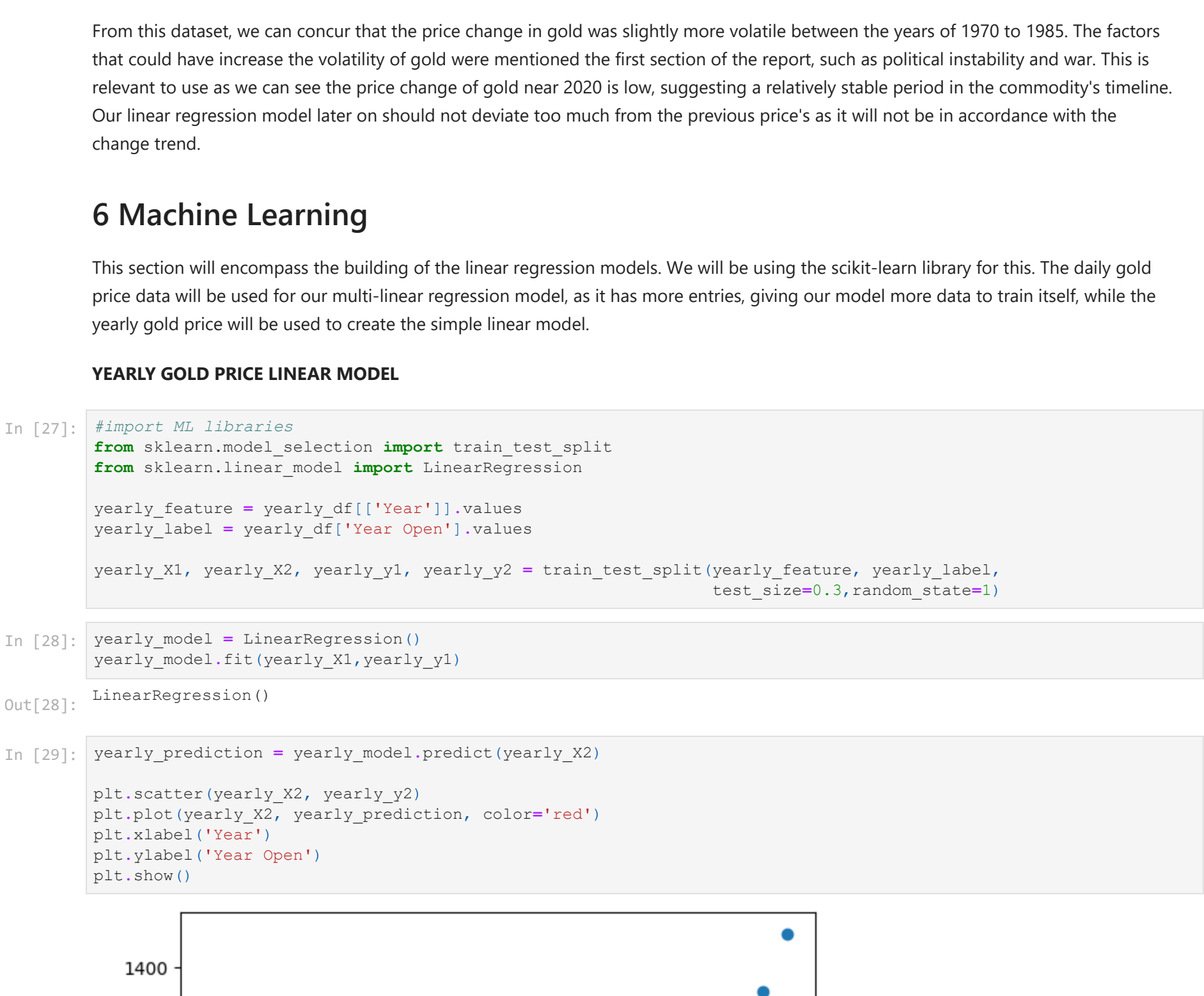
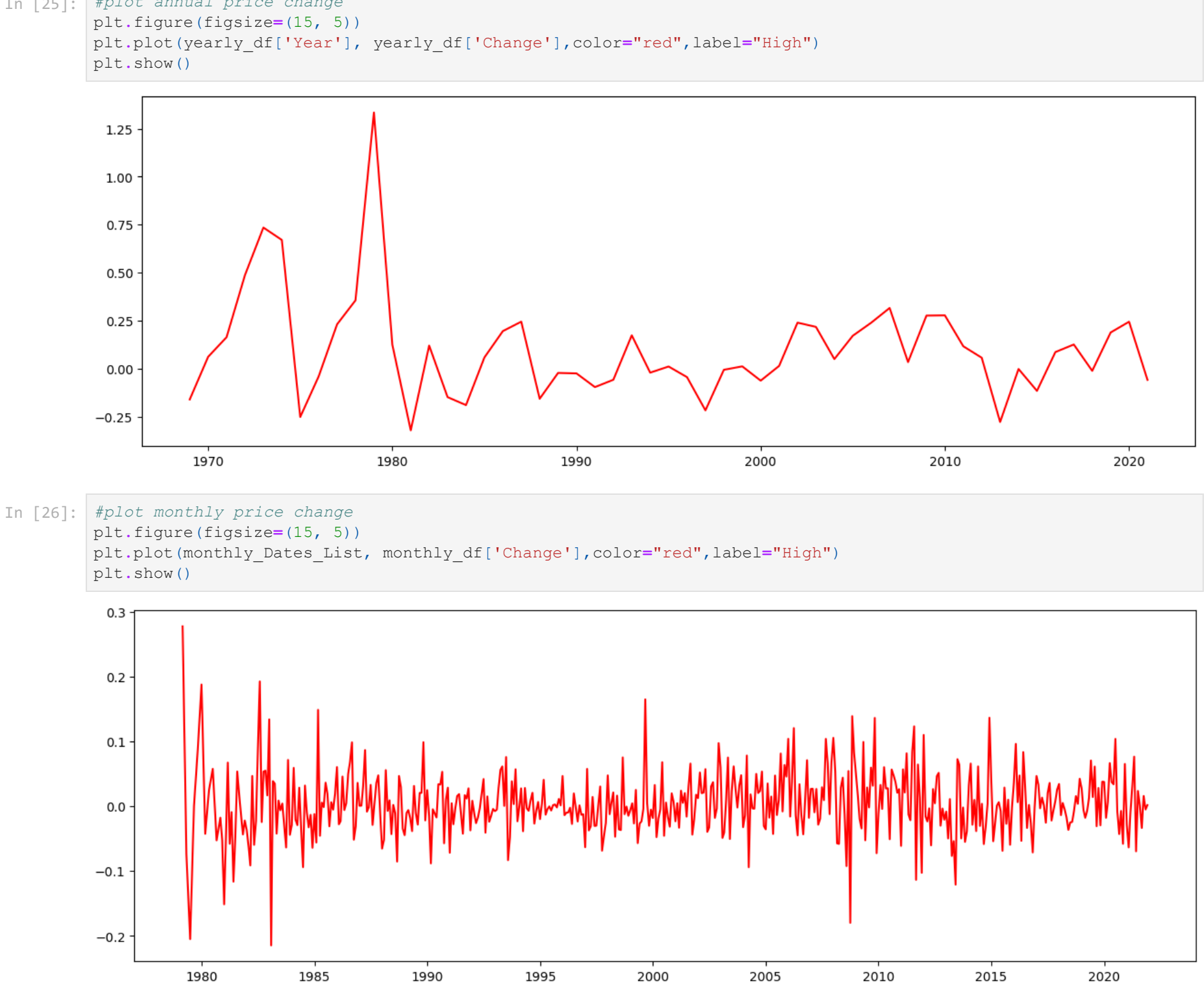
The monthly volume describes a fairly similar pattern.

5 Visualisation

```
In [23]: #plot average gold yearly closing price
fig, axes = plt.subplots(3, figsize=(7, 10))
fig.tight_layout()

axes[0].plot(yearly_df['Yearly'], yearly_df['Average'])
axes[0].title.set_text('Yearly Closing Price')
axes[1].plot(monthly_Dates_List, daily_df['High'], color='red', label='High')
axes[1].title.set_text('Yearly Highs')
axes[2].plot(daily_Dates
```


The green line shows the Low price of the gold for the day, whilst the red line in the back, highlights the high of the day. The close correlation between the daily high and lows can imply a tight bid-ask spread[8]. Which means gold was actively traded throughout the decades

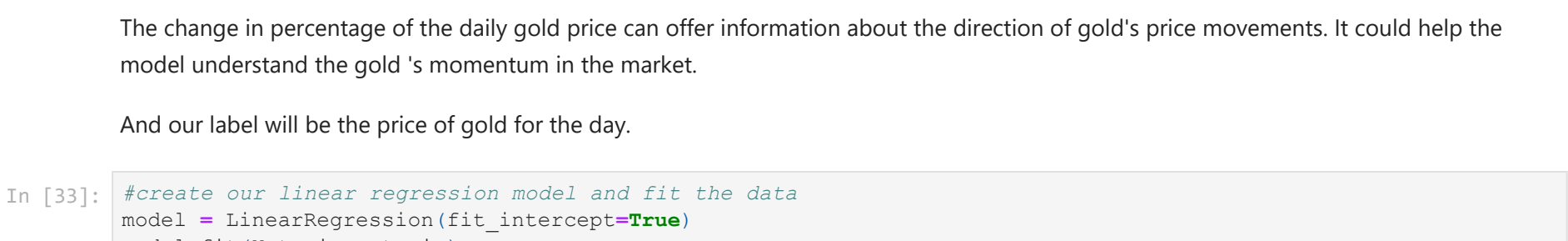
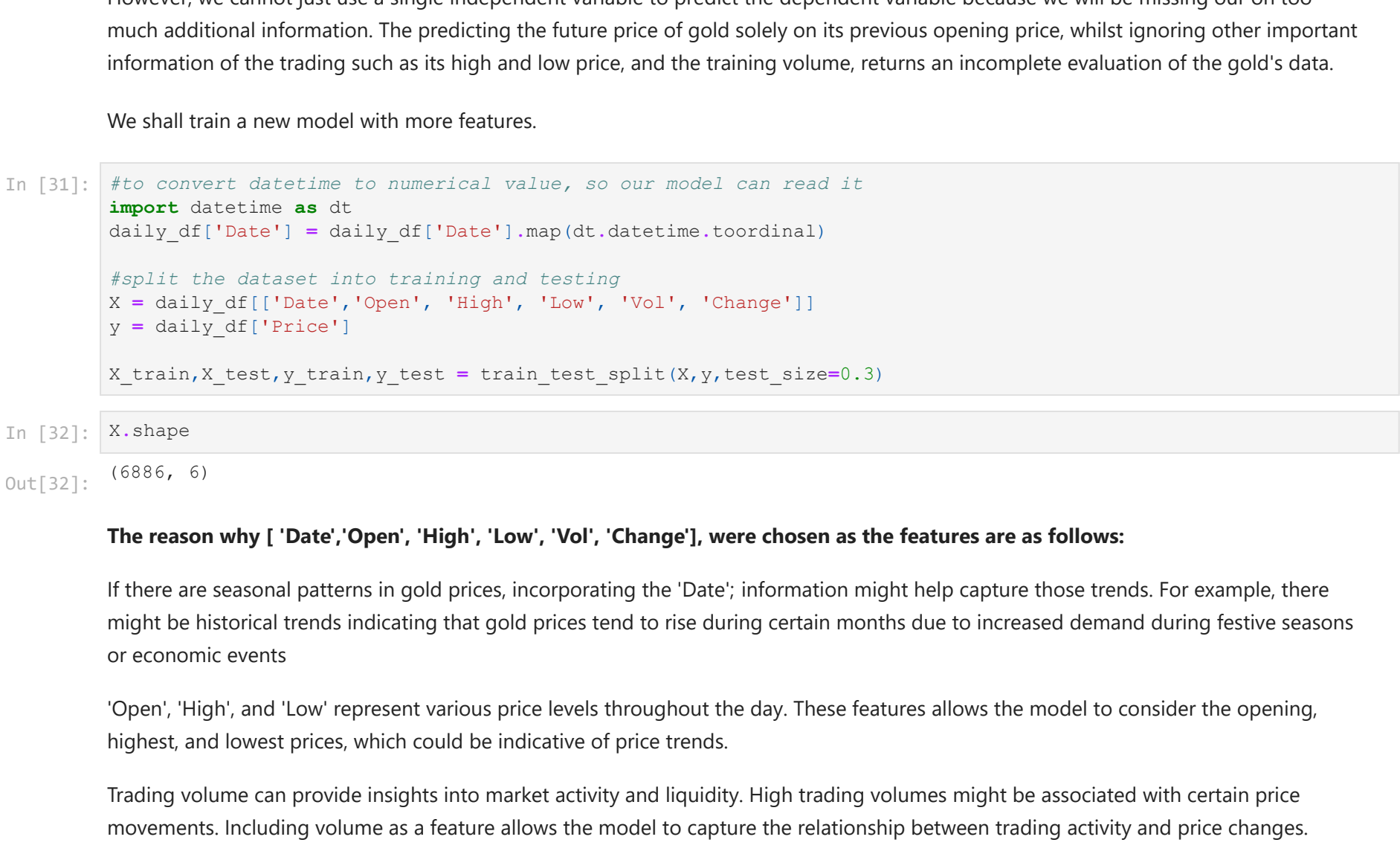


From this dataset, we can concur that the price change in gold was slightly more volatile between the years of 1970 to 1985. The factors that could have increased the volatility of gold were mentioned the firsty more of the report, such as political instability and war. This is relevant to use as we can see the price change of gold near 2020 is low, suggesting a relatively stable period in the commodity's timeline. Our linear regression model later on should not deviate too much from the previous price's as it will not be in accordance with the change trend.

6 Machine Learning

This section will encompass the building of the linear regression models. We will be using the scikit-learn library for this. The daily gold price data will be used for our multi-linear regression model as it has more entries, giving our model more data to train itself, while the yearly gold price will be used to create the simple linear model.

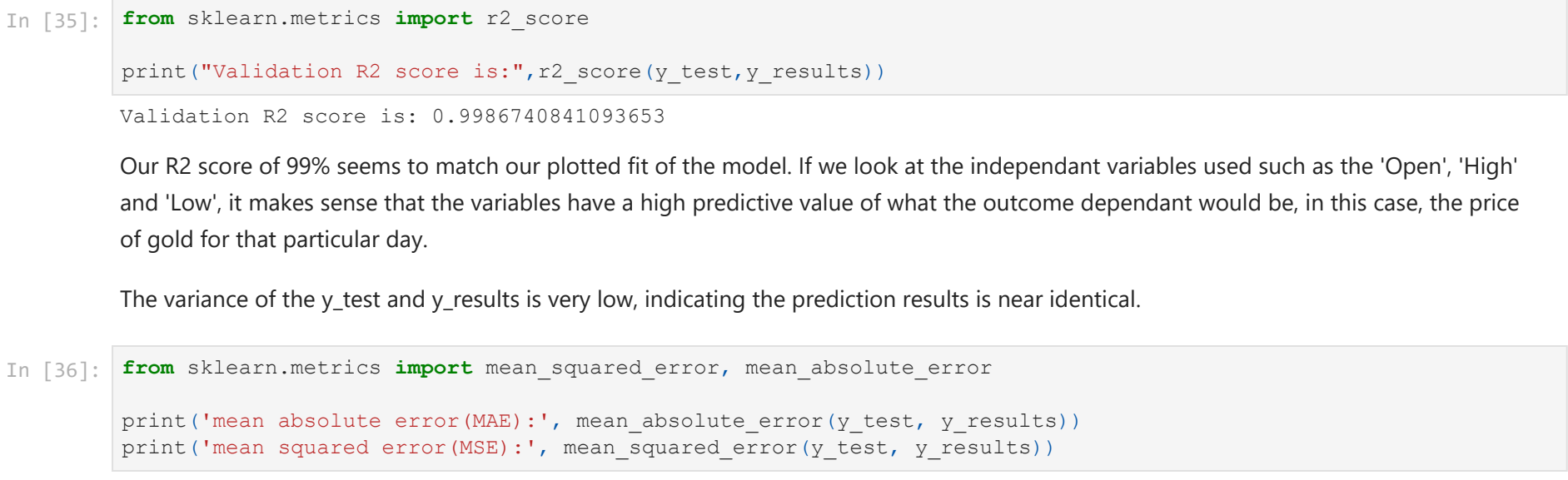
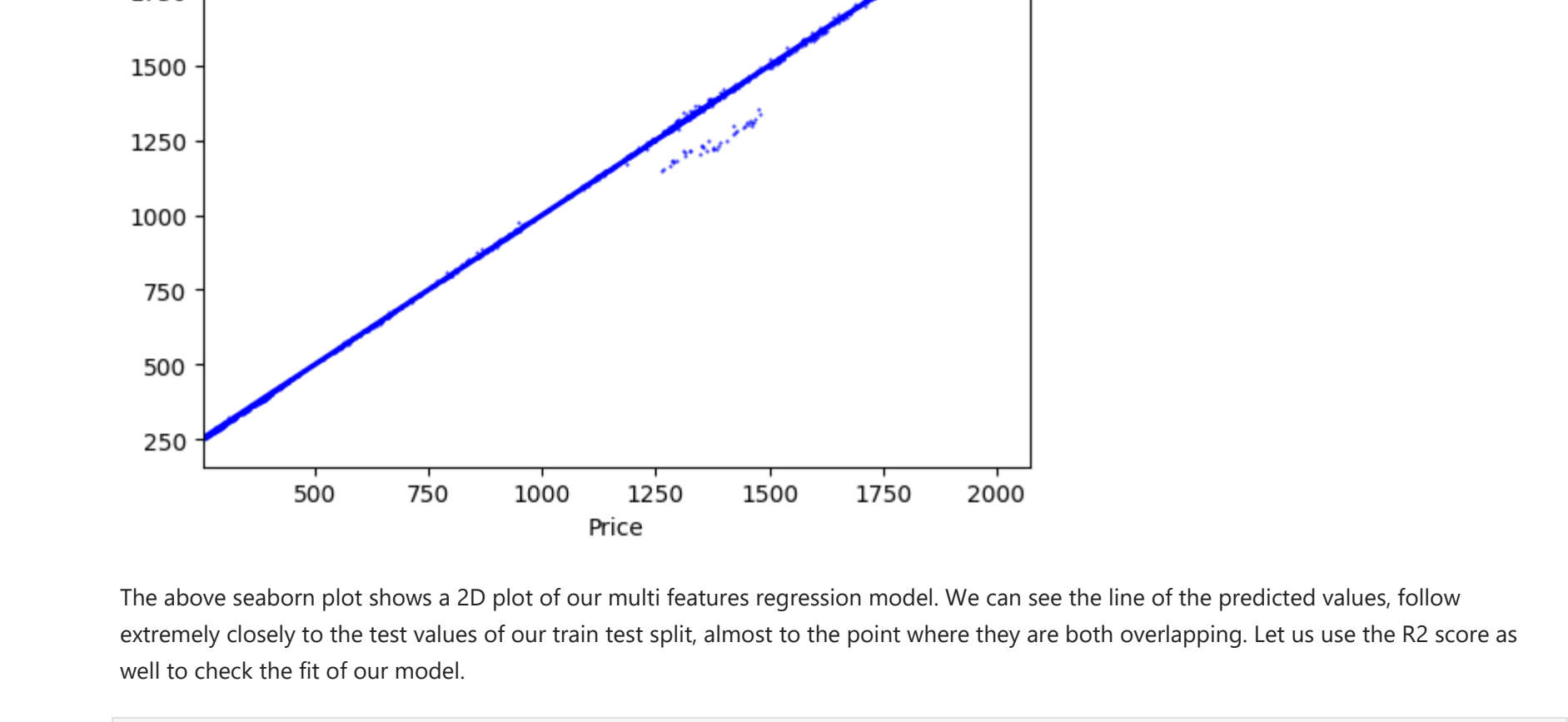
YEARLY GOLD PRICE LINEAR MODEL



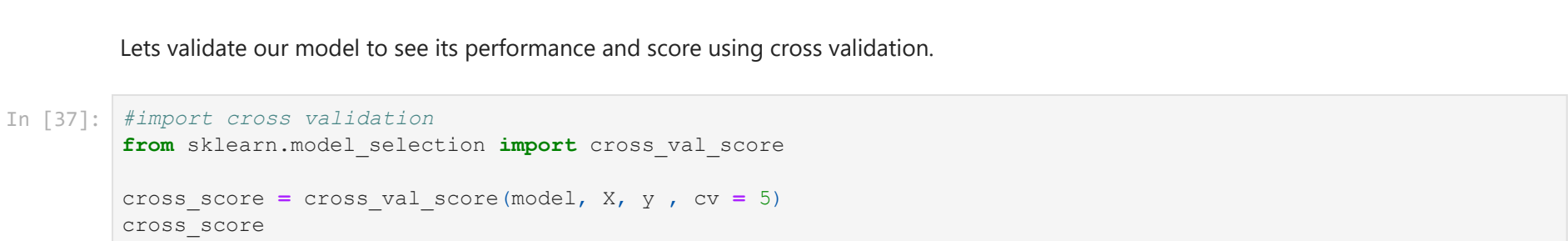
DAILY GOLD PRICE MULTI LINEAR MODEL

However, we cannot just use a single independent variable to predict the dependent variable because we will be missing our on too much additional information. The predicting the future price of gold solely on its previous opening price, whilst ignoring other important information of the trading such as its high and low price, and the trading volume, returns an incomplete evaluation of the gold's data.

We shall train a new model with more features.

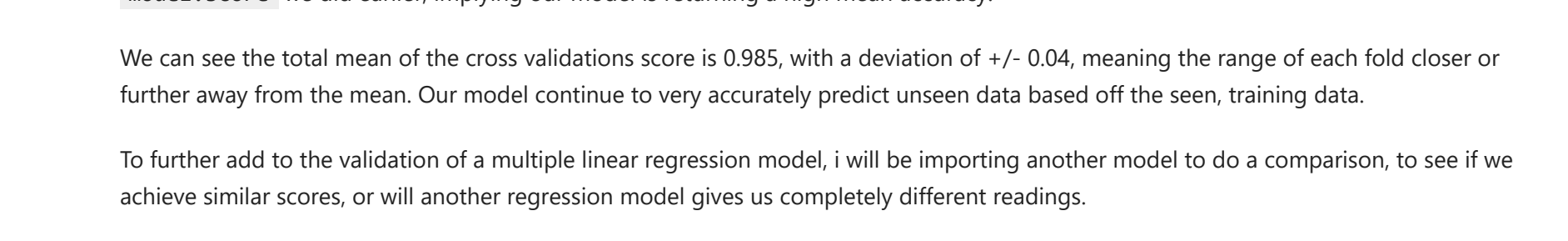


The above seaborn plot shows a 2D plot of our multi features regression model. We can see the line of the predicted values, follow extremely closely to the test values of our train test split, almost to the point where they are both overlapping. Let us use the R2 score as well to check the fit of our model.



Our R2 score of 99% seems to match our plotted fit of the model. If we look at the independent variables used such as the 'Open', 'High' and 'Low', it makes sense that the variables have a high predictive value of what the outcome dependent would be, in this case, the price of gold for that particular day.

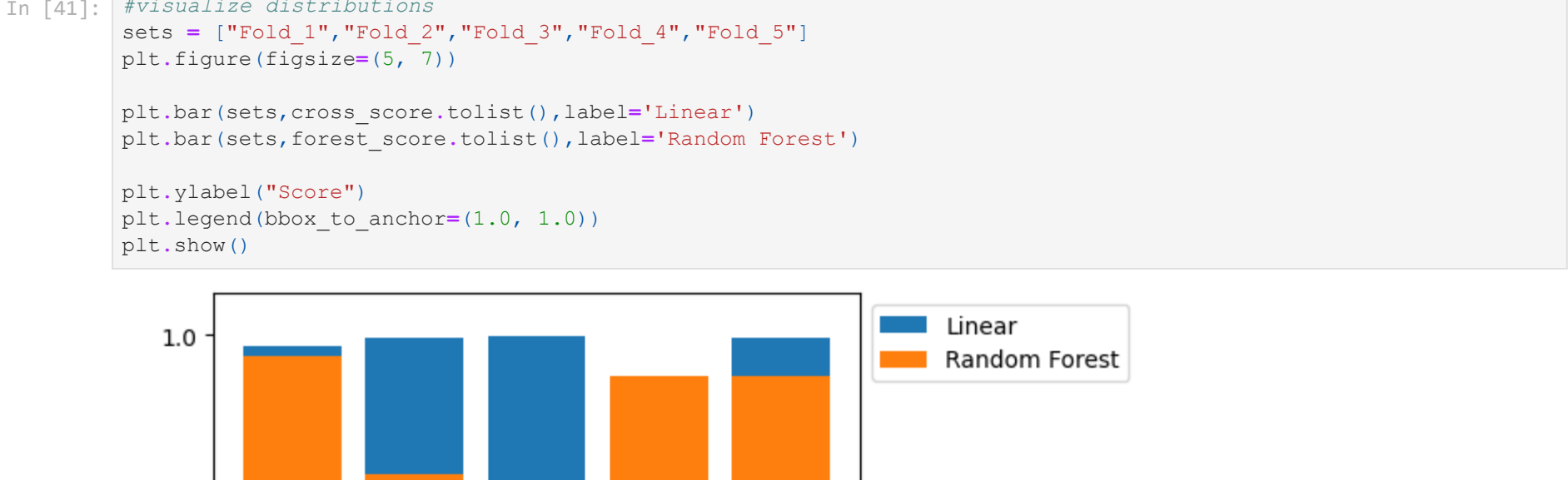
The variance of the y_test and y_results is very low, indicating the prediction results is near identical.



The mean absolute error (MAE) is defined as the sum of absolute/positive errors of all values. We can see the absolute error of our test data and the predicted values being off about 6, whilst the mean squared error is the sum of squared differences, the lower the better.

7 Validation

Lets validate our model to see its performance and score using cross validation.

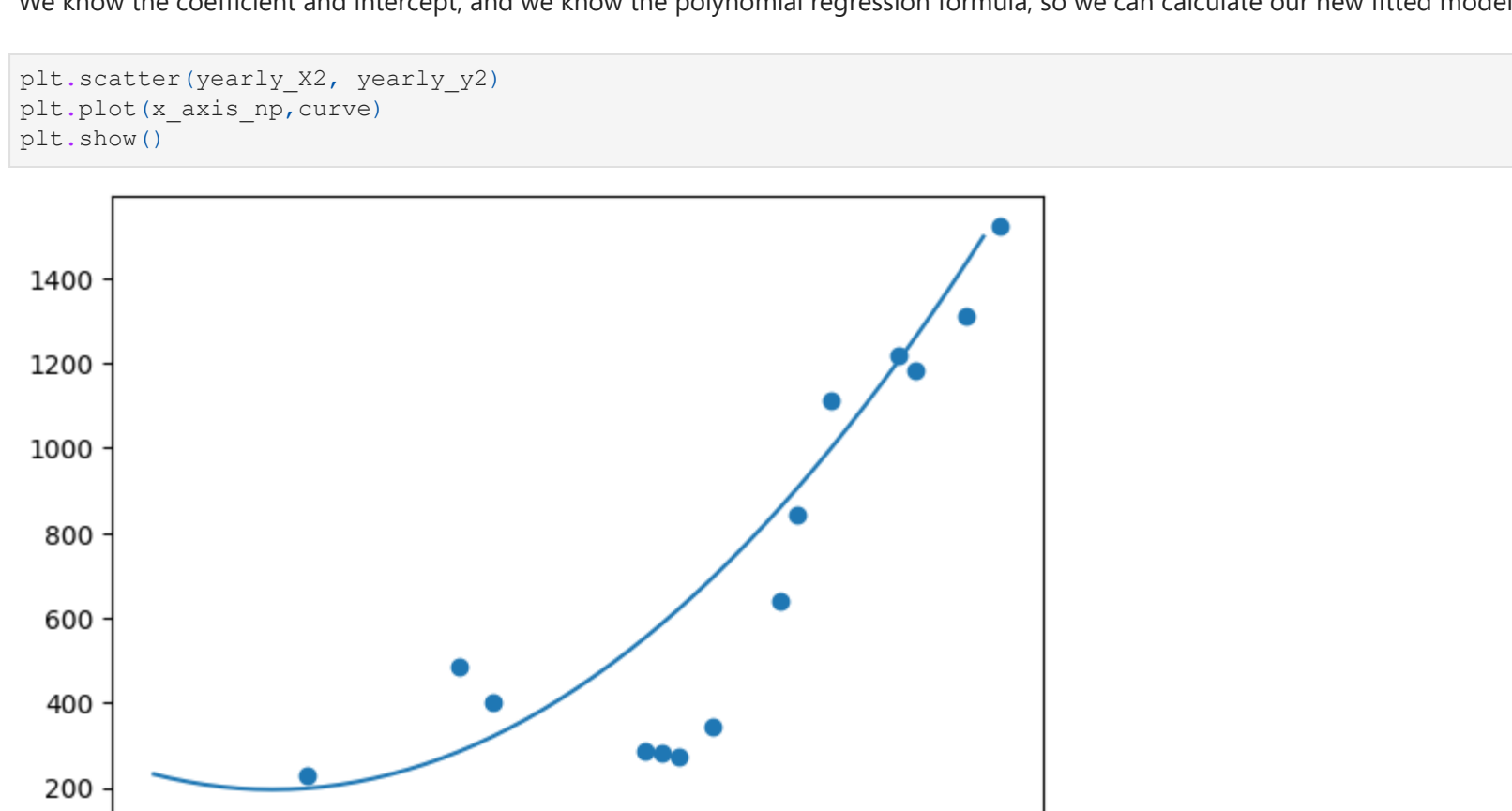
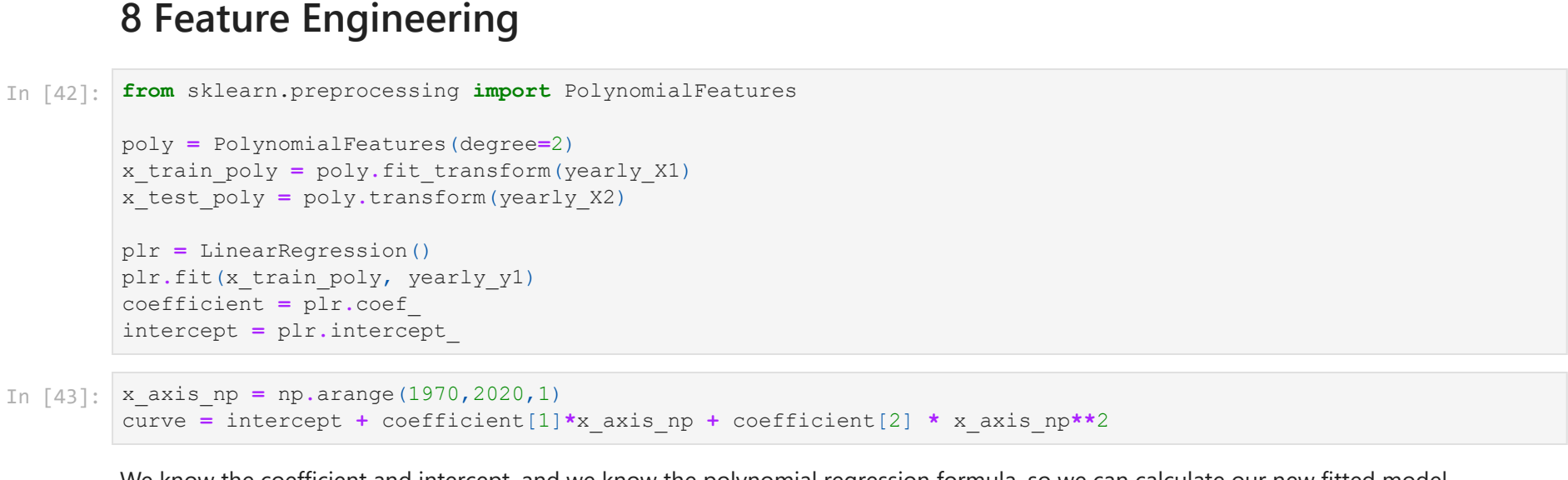


The array return by cross_val_score is the score of the estimator for each one of the five sets of the cross validation[9]. We can see the score of each of the five sets, is very close to 1, which is the maximum score attainable. This score closely matches the initial model score, we did earlier, implying our model is returning a high mean accuracy.

We can see the total mean of the cross validations score is 0.985, with a deviation of +/- 0.04, meaning the range of each fold closer or further away from the mean. Our model continue to very accurately predict unseen data based off the seen, training data.

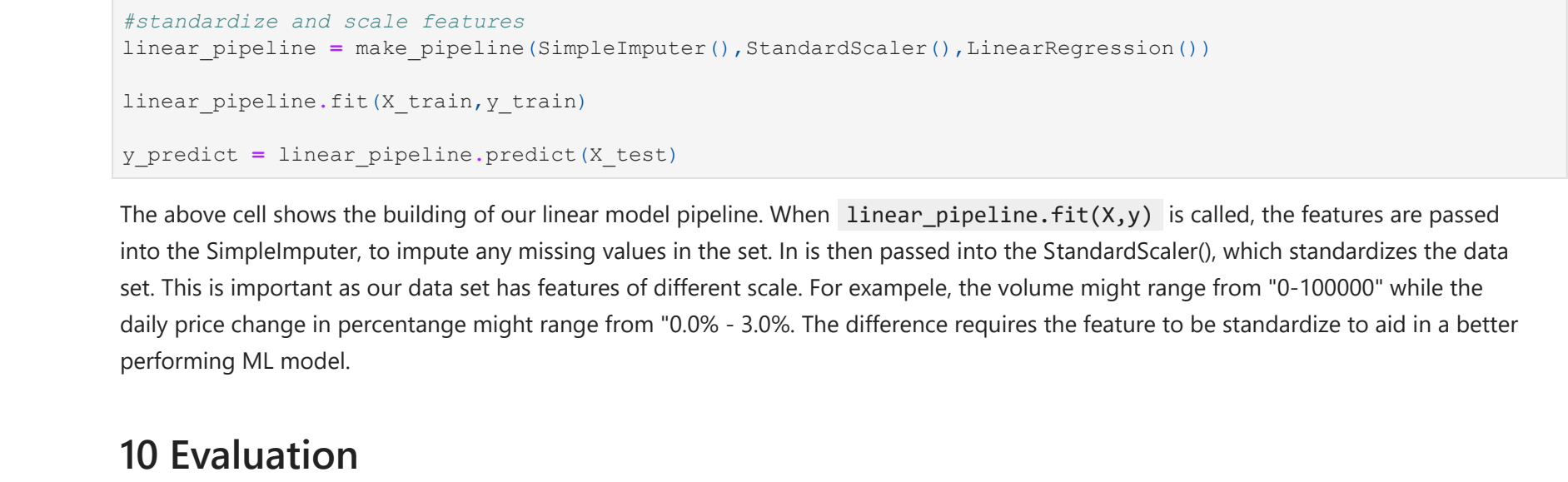
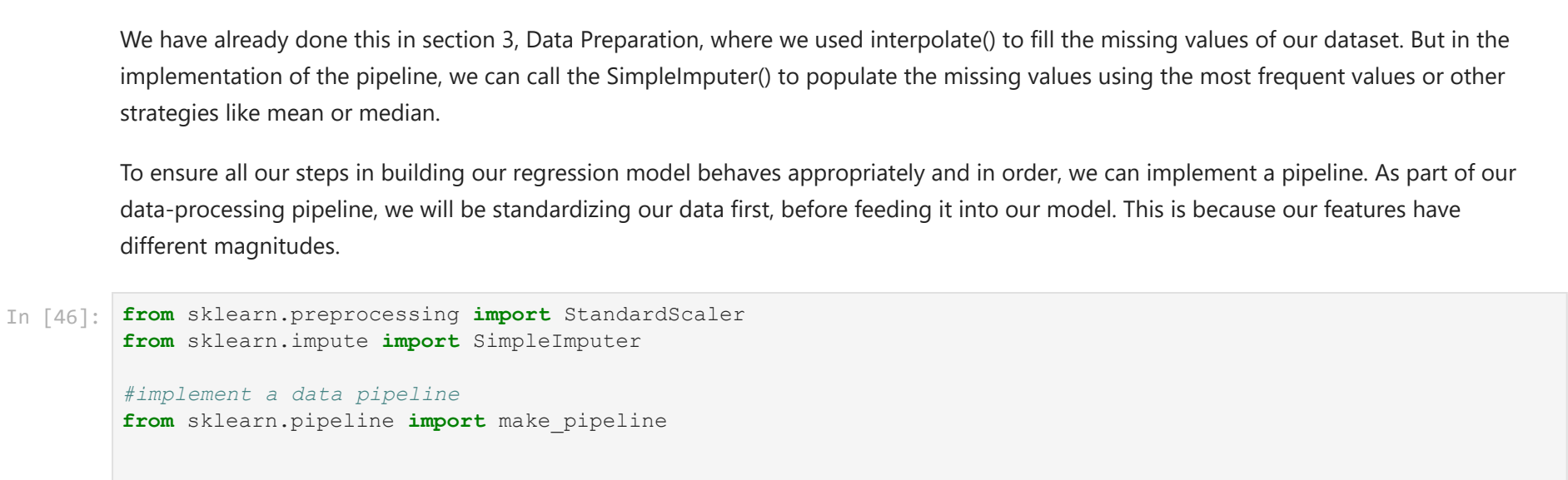
To further add to the validation of a multiple linear regression model, i will be importing another model to do a comparison, to see if we achieve similar scores, or will another regression model gives us completely different readings.

The model I have chosen is the Random Forest Regression Model. The reason is because this model uses averaging to improve the predictive accuracy and control over-fitting[10].

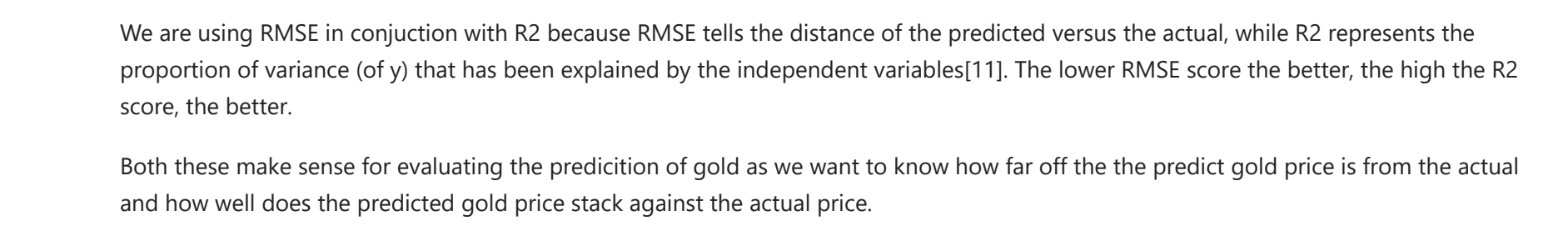


Firstly, the forest cross validation score mean was 0.860 with a deviation of 0.29. Our linear model gave a better validation of 0.995 and 0.04 respectively. The score difference can be seen above prior to building the pipeline, measures the difference between the predicted versus the actual values, squares it, then returns the mean of all the samples in that set.

8 Feature Engineering



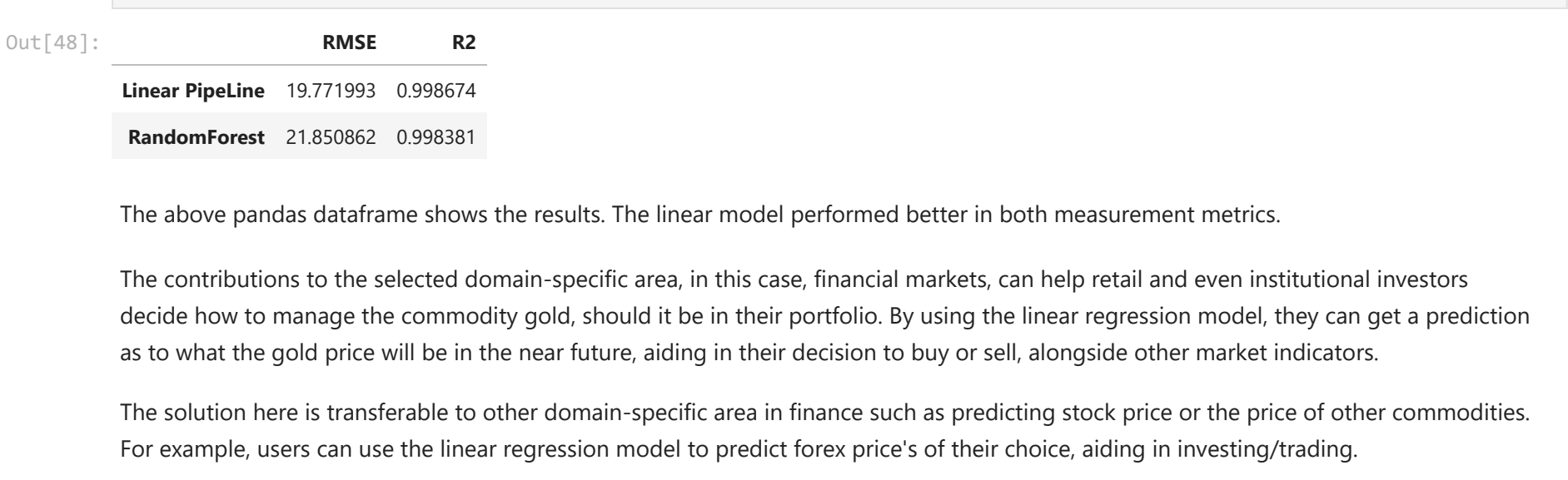
We can see after adding polynomial degree of 2, the yearly gold price prediction model fits a much better line to the training data, compared to without it.



An increase R2 score compared to the default model. The feature engineering of polynomial degrees has improved our model.

Our models has already benefited from one of the feature engineering techniques described in the course, imputation of missing values. We have already done this in section 3. Data Preparation, where we used interpolate() to fill the missing values of our dataset. But in the implementation of the pipeline, we can call the SimpleImputer() to populate the missing values using the most frequent values or other strategies like mean or median.

To ensure all our steps in building our regression model behaves appropriately and in order, we can implement a pipeline. As part of our data-processing pipeline, we will be standardizing our data first, before feeding it into our model. This is because our features have different magnitudes.



The above cell shows the building of our linear model pipeline. When linear_pipeline.fit(X,y) is called, the features are passed into the SimpleImputer, to impute any missing values in the set. In is then passed into the StandardScaler(), which standardizes the data set. This is important as our data set has features of different scale. For example, the volume might range from "0-100000" while the daily price change in percentage might range from "0.0% - 3.0%". The difference requires the feature to be standardized to aid in a better performing ML model.

10 Evaluation

To close of this final section of the report, we shall evaluate the simple linear regression model and the multi linear model built through the data pipeline against the validation random forest model.

To evaluate, we will be measuring the R-Square(R2) AND Root Mean Square Error(RMSE) of the model, which is the root of the mean square error. The mean square error, as seen above prior to building the pipeline, measures the difference between the predicted versus the actual values, squares it, then returns the mean of all the samples in that set.

We are using RMSE in conjunction with R2 because RMSE tells the distance of the predicted versus the actual, while R2 represents the proportion of variance (of y) that has been explained by the independent variables[11]. The lower RMSE score the better, the higher the R2 score, the better.

Both these make sense for evaluating the prediction of gold as we want to know how far off the the predicted gold price is from the actual and how well does the predicted gold price stack against the actual price.



The above pandas dataframe shows the results. The linear model performed better in both measurement metrics.

The contributions to the selected domain-specific area in this case, financial markets, can help retail and even institutional investors decide how to manage the commodity gold, should it be in their portfolio. By using the linear regression model, they can get a prediction as to what the gold price will be in the near future, aiding in their decision to buy or sell, alongside other market indicators.

The solution here is transferable to other domain-specific area in finance such as predicting stock price or the price of other commodities. For example, we can use the linear regression model to predict forex price's of their choice, aiding in investing/trading.

References

- [1] World Gold Council. (n.d.). Money and Gold. [online] Available at: <https://www.gold.org/history-gold/gold-as-currency>
- [2] ETF Database. (n.d.). Gold ETF List. [online] Available at: <https://etfdb.com/etfs/commodity/gold/>
- [3] Maxwell, T. (2022). Why You Should Buy Gold during Inflation. [online] www.cbsnews.com. Available at: <https://www.cbsnews.com/news/why-you-should-buy-gold-during-inflation/>
- [4] Dataset : https://www.kaggle.com/datasets/nward7/gold-historical-datasets?select=Gold_Monthly.csv
- [5] Investing.com (2023). Investing.com - Stock Market Quotes & Financial News. [online] Investing.com. Available at: <https://www.investing.com/>.
- [6] Chen, J. (2019). Bull Market Definition. [online] Investopedia. Available at: <https://www.investopedia.com/terms/b/bullmarket.asp>.
- [7] Shumsky, T. (2013). Gold Falls 28% In 2013, Ends 12-Year Bull Run. Wall Street Journal. [online] 31 Dec. Available at: <https://www.wsj.com/articles/SB10001424052702304591604579292321014055380>
- [8] Investopedia. (n.d.). What Is a Bid-Ask Spread, and How Does It Work in Trading? [online] Available at: <https://www.investopedia.com/terms/b/bid-ask-spread.asp?~:text=In%20financial%20markets%2C%20a%20bid>.
- [9] Scikit-learn.org. (2019). sklearn.model_selection.cross_val_score - scikit-learn 0.22 documentation. [online] Available at: https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.cross_val_score.html.
- [10] scikit-learn (2018). 3.2.4.3.2. sklearn.ensemble.RandomForestRegressor - scikit-learn 0.20.3 documentation. [online] Scikit-learn.org. Available at: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html>.
- [11]scikit-learn. (n.d.). 3.3. Metrics and scoring: quantifying the quality of predictions. [online] Available at: https://scikit-learn.org/stable/modules/model_evaluation.html#2-score