# Oil & Gold Price Forecasting
# using VAR and Granger Causality

*By* Group 6[*]

Jaihind (19BM6JP43), Deo Ashish Samanta (19BM6JP44)
Sarnesh Kalundia (19BM6JP45), Vineet Kumar(19BM6JP46)
Soham Mullick (19BM6JP47), Koustav Mitra (19BM6JP48)
Rajat Gautam (19BM6JP49), Kowshik Moyya (19BM6JP50), and
Chetan Mehatre (19BM6JP60)

*Oil & Gold price forecasting has become an important problem for academics as well as other stakeholders of the market. In this project we revisit the research question of significant relationship between Indian economy and historical oil price. We aim to price forecast (in INR) and assess the interdependence & causality between IIP index, oil import & export by India, and precious metal price. We are using a vector auto-regressive model to forecast and extract relevant inferences. We are using a 12 year rich monthly data sourced from EPWRF database and evaluated our model on 15 step ahead forecast using standard metrics viz. RMSE, MAPE etc.*

## I.  Introduction

Energy resources like oil occupy a crucial position in the modern economics as they are extensively used in the provisions of goods and services. That is why volatility in their prices has been a major concern for modern economies especially since World War II. Results from existing literature on the relationship between oil price and macro economy are mixed.

Gold is used as a risk management benchmark tool in hedging and diversifying commodity portfolios. Investors in advanced and emerging markets often switch between oil and gold or combine them to diversify their portfolios

Forecasting Gold and Oil has garnered major attention from academics, investors and Government agencies alike. These two products are known for their substantial influence on the global economy. However, what is relationship between oil price and macro economy in Indian context with updated data and methodology, What generally happens to Gold prices when Oil takes a plunge? In this project we use a vector auto-regressive model to do forecasting of future Gold & Oil prices.

## II.   Dataset

The dataset is downloaded from http://www.epwrfits.in, We visualize all the time series under consideration to get a preliminary understanding of their nature. From the plots we can see that most of the time series are having an increasing trend and thus are non stationary. We also see the presence of seasonality in the Mining Index. Also we observe that many of the series seem to have a positive correlation between pairs of such time series variables like Electricity and Coke Fuel Index, Oil Import and Export and Oil Price and WPI for mineral oils

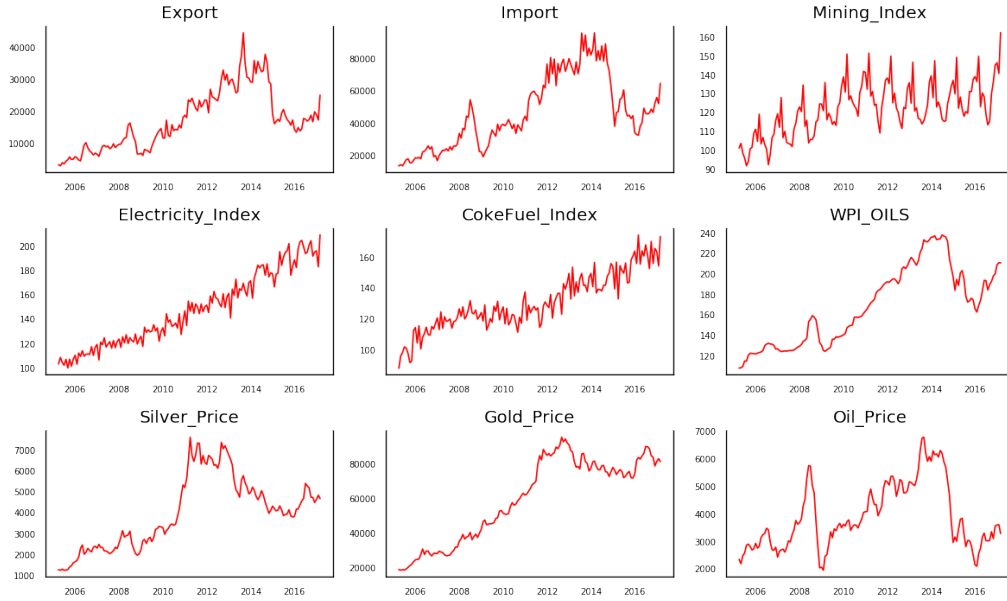The raw dataset is plotted and shown in Figure: 1



Figure 1. : Raw Time Series Data

### A.   Stationarity

Stationarity can be defined as a constant mean, variance or autocovariances for each given lag. There are several reasons why the concept of non-stationarity is important and why it is essential that variables that are non-stationary should be treated differently from variables that are stationary. These reasons are given as follows:

- The stationarity or otherwise of a series can strongly influence its behavior or properties. To offer one illustration, the word 'shock' usually used to denote a change or an unexpected change in a variable or perhaps simply the value of the error term during a particular time period. For a stationary
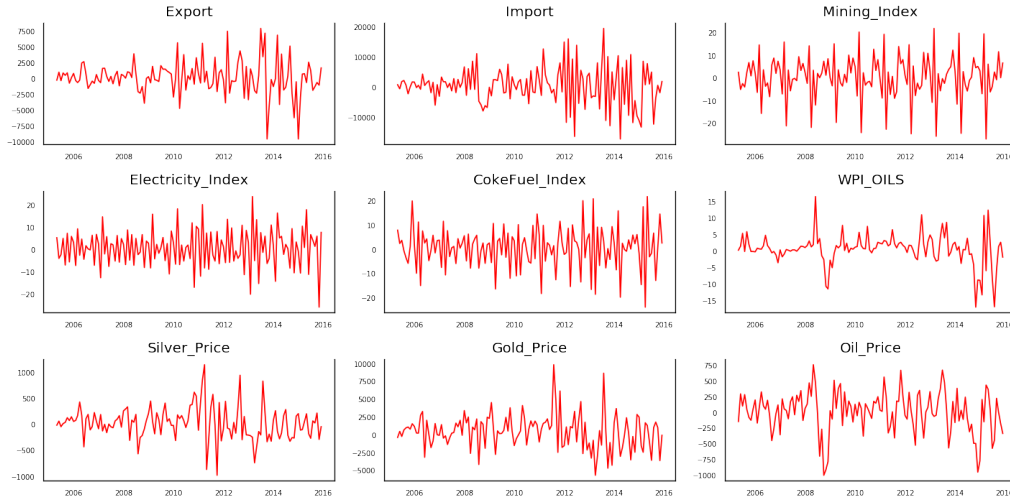
Figure 2. : Stationary Time Series after First Differencing

series, 'shocks' to the system will gradually die away. This can be contrasted with the case of non-stationary data, where the persistence of shocks will always be infinite.

- The use of non-stationarity data can lead to spurious regressions. If two stationary variables are generated as independent random series, when one of those variables is regressed on the other, the t-ratio on the slope coefficient would be expected not to be significantly different from zero, and the value of $R^2$ would be expected to be very low. So, if standard regression techniques are applied to non-stationary data, the end result could be a regression that 'looks' good under standard measures, but which is really valueless. Such a model would be termed a 'spurious' regression.

First order difference makes the series stationary. The same is shown in Figure 2.

## III.   Causation Analysis

As we have included multiple time series we want to see how many of theses time series actually aid us in improving predictive performance of models.

For this, we use granger causality test — A time series X is said to Granger-cause Y if it can be shown, usually through a series of t-tests and F-tests on lagged values of X (and with lagged values of Y also included), that those X values provide statistically significant information about future values of Y.

TESTING GRANGER CAUSALITY IDEA.  — The test is pretty simple we include another time series that we want to assess and test for corresponding coefficients

of regression. If the p values is small we reject the null hypothesis that is $X$, do not cause other series, $Y$.

More formally, Let y and x be stationary time series. To test the null hypothesis that x does not Granger-cause y, one first finds the proper lagged values of y to include in a univariate autoregression of y:

$$y_t = a_0 + a_1 y_{t-1} + a_2 y_{t-2} + \cdots + a_m y_{t-m} + \text{error}_t.$$

Next, the auto-regression is augmented by including lagged values of x:

$$y_t = a_0 + a_1 y_{t-1} + a_2 y_{t-2} + \cdots + a_m y_{t-m} + b_p x_{t-p} + \cdots + b_q x_{t-q} + \text{error}_t.$$

One retains in this regression all lagged values of x that are individually significant according to their t-statistics, provided that collectively they add explanatory power to the regression according to an F-test (whose null hypothesis is no explanatory power jointly added by the x's). In the notation of the above augmented regression, p is the shortest, and q is the longest, lag length for which the lagged value of x is significant.

The null hypothesis that x does not Granger-cause y is accepted if and only if no lagged values of x are retained in the regression.

we must also run this model in reverse to verify that $Y$ does not provide information about future value of $X$.If we find that this is the case, it is likely that there is some exogenous variable, $Z$, which needs to be controlled or could be a better candidate for Granger causation.

| | Export_x | Import_x | Mining_Index_x | Electricity_Index_x | CokeFuel_Index_x | WPI_OILS_x | Silver_Price_x | Gold_Price_x | Oil_Price_x |
|---|---|---|---|---|---|---|---|---|---|
| Export_y | 1.000 | 0.000 | 0.001 | 0.000 | 0.228 | 0.000 | 0.000 | 0.000 | 0.000 |
| Import_y | 0.000 | 1.000 | 0.008 | 0.153 | 0.075 | 0.000 | 0.000 | 0.001 | 0.000 |
| Mining_Index_y | 0.000 | 0.000 | 1.000 | 0.000 | 0.000 | 0.008 | 0.005 | 0.000 | 0.036 |
| Electricity_Index_y | 0.013 | 0.057 | 0.000 | 1.000 | 0.000 | 0.001 | 0.050 | 0.001 | 0.168 |
| CokeFuel_Index_y | 0.001 | 0.002 | 0.028 | 0.000 | 1.000 | 0.000 | 0.015 | 0.000 | 0.032 |
| WPI_OILS_y | 0.013 | 0.000 | 0.020 | 0.135 | 0.017 | 1.000 | 0.000 | 0.006 | 0.000 |
| Silver_Price_y | 0.242 | 0.151 | 0.080 | 0.633 | 0.277 | 0.151 | 1.000 | 0.342 | 0.256 |
| Gold_Price_y | 0.294 | 0.268 | 0.120 | 0.802 | 0.292 | 0.292 | 0.033 | 1.000 | 0.232 |
| Oil_Price_y | 0.163 | 0.082 | 0.111 | 0.346 | 0.275 | 0.206 | 0.031 | 0.103 | 1.000 |

Figure 3. : Granger Causation Matrix

In figure 4 we can see the table of p values; columns are predictors and rows are a response. Most of the variables are interchangeably causing each other. This justifies the VAR modeling approach.

## IV.   Co-integration and Johansen Test

Regression errors are a linear combination of two unit root series ($X$ and $Y$) of the following equation:
$$\varepsilon = Y_t - \alpha - \beta X_t$$
If $X_t$ and $X_t$ have a unit root, this usually means $\varepsilon_t$ also has a unit root. This is the cause of a spurious regression problem.

However, there are some cases when the unit root in $Y$ and $X$ cancel each other out and $\varepsilon_t$ does not have any unit root. In this special case, the spurious regression problem vanishes. This phenomenon is referred to as cointegration, which measures the equilibrium spread between two non-stationary time series measured by constant terms $\alpha$. If the error term is stationary then we can say that cointegration and an equilibrium spread exist.

**Definition.** The two $I(d)$ processes (Time Series) $X_t$ and $Y_t$ are said to have cointegration, if there exist a linear combination $aX_t + bY_t$ which is an $I(m)$ process for some $m < d$.

### A.   Johansen Test

Consider the joint hypothesis

$$H_0 : r = 0 \text{ vs } H_1 : 0 < r \le g$$

$$H_0 : r = 1 \text{ vs } H_1 : 1 < r \le g$$
$$H_0 : r = 2 \text{ vs } H_1 : 2 < r \le g$$
$$\ldots \qquad \ldots \qquad \ldots$$
$$H_0 : r = g - 1 \text{ vs } H_1 : r = g$$

**Test Statistics:**

$$\lambda_{\text{trace}}(r) = -T \sum_{i=r+1}^{g} \ln\left(1 - \hat{\lambda}_i\right) \text{ and } \lambda_{\max}(r, r+1) = -T \ln\left(1 - \hat{\lambda}_{r+1}\right)$$

where, $\widehat{\lambda}_i$ is the estimated value for the $i$th ordered eigenvalue from the $\Pi$ matrix. To use Johansen's method, we need to turn the VAR of the form

$$y_t = \beta_1 y_{t-1} + \beta_2 y_{t-2} + \ldots + \beta_k y_{t-k} + u_t$$

into a VECM, which can be written as

$$\Delta y_t = \Pi y_{t-1} + \Gamma_1 \Delta y_{t-1} + \Gamma_2 \Delta y_{t-2} + \ldots + \Gamma_{k-1} \Delta y_{t-(k-1)} + u_t$$

```
Name    :: Test Stat > C(95%)        => Signif
  ----------------------------------------
Export :: 360.23    > 179.5199    =>    True
Import :: 253.07    > 143.6691    =>    True
Mining_Index :: 174.24    > 111.7797    =>    True
Electricity_Index :: 118.99    > 83.9383    =>    True
CokeFuel_Index :: 75.22    > 60.0627    =>    True
WPI_OILS :: 46.43    > 40.1749    =>    True
Silver_Price :: 27.3    > 24.2761    =>    True
Gold_Price :: 13.19    > 12.3212    =>    True
Oil_Price :: 5.86    > 4.1296    =>    True
```

Figure 4. : Johansen Test corresponding to each variable

*Upon performing the test we could reject the Null Hypothesis, for $r = 0$, hence* $\Pi$ *has rank* $0$. **There are no cointegration variables**.

## V.   Vector Auto Regressor Model

As we can see in previous section, there is no co-integration in any of our variables. Now we can conduct Vector auto-regressor or VAR, which is essentially a multivariate forecasting algorithm used when two or more time-series influence one another.

In the VAR model, each variable is modeled as a linear combination of past values of itself and the past values of other variables in the system.

More formally, VAR models are characterized by their order, which refers to the number of earlier time periods the model will use. A lag is the value of a variable in a previous time period. So in general a pth-order VAR refers to a VAR model which includes lags for the last p time periods. A pth-order VAR is denoted "VAR(p)" and sometimes called "a VAR with p lags". A pth-order VAR model is written as

$$y_t = c + A_1 y_{t-1} + A_2 y_{t-2} + \cdots + A_p y_{t-p} + e_t$$

The variables of the form $y_{t-i}$ indicate that variable's value $i$ time periods earlier and are called the "ith lag" of $y_t$. The variable $c$ is a $k-$vector of constants serving as the intercept of the model. $A_i$ is a time-invariant $(k \times k)$ matrix and $e_t$ is a $k-$vector of error terms.

The process of choosing the maximum lag $p$ in the VAR model requires special attention because inference is dependent on correctness of the selected lag order.

As you increase the number of time series in the model, the system of equations become larger. What we will attempt to do is fit the VAR model on the training set and then use the fitted model to forecast the next 15 observations. These forecasts will be compared to the observations in the testing dataset. We have taken the maximum lag of 5 to identify the required lags in the VAR model.

## VI.  Forecasting

After deciding on the model we forecasted the price of Oil as well as Gold using the VAR model on the above mentioned test data. To evaluate the forecast in a more comprehensive manner, we will use mean absolute error, mean squared error, and root mean squared error. What are the difference between these three accuracy metrics:

- Mean Absolute Error (MAE) measures the average magnitude of the errors in a set of predictions, without considering their direction. It's the average over the test sample of the absolute differences between prediction and actual observation where all individual differences have equal weight

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|$$

- Root Mean Squared Error (RMSE) is a quadratic scoring rule that also measures the average magnitude of the error. It's the squared root of the average of squared differences between prediction and actual observation.,

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2}$$

Taking the square root of the average squared errors has some interesting implications for RSME. Since the errors are squared before they are averaged, the RMSE gives a relatively high weight to large errors. This means the RMSE sould be more useful when large errors are particularly undesirable.

- Mean absolute percentage error (MAPE), also known as mean absolute percentage deviation (MAPD), is a measure of prediction accuracy of a forecasting method in statistics

$$M = \frac{1}{n} \sum_{t=1}^{n} \left| \frac{A_t - F_t}{A_t} \right|,$$

where where $A_t$ is the actual value and $F_t$ is the forecast value.

Higher value of RMSE in gold is due to the higher magnitude of the original values, not because of forecast. Which is evident from the 2% MAPE value in Table 1.

We have also plotted our forecast value and overlayed it with original series in Figure 5.

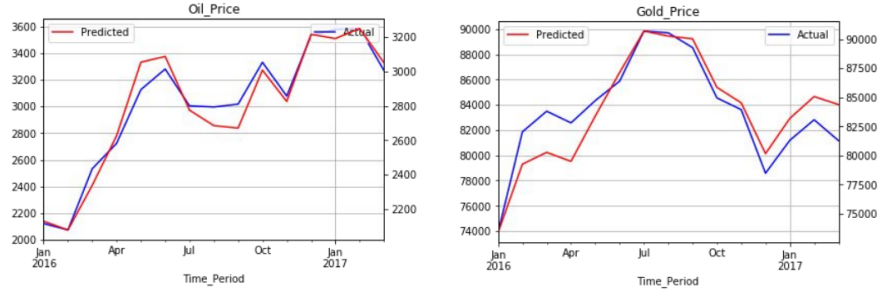|  | Gold Price | Oil Price |
|---|---|---|
| MAPE | 2.067% | 6.864% |
| MAE | 1714.430 | 218.703 |
| RMSE | 1947.669 | 249.501 |

Table 1—: Forecast Performance



Figure 5. : Time Series Forecast

## VII.    Discussion and Insights

In the current study, we tried to identify and explain any causalities among gold and oil prices with other economic factors like exports and imports, WPI, electricity, and mining index values and mining index values for India, using Granger causality theory. Causal connectivity implies a connection that lies in the enhancement of forecasting between two time-series. That being said and according to our knowledge, the lack of similar research on non-market coupled countries makes the current work unique and the outcomes that arise should be taken under consideration

The model that was presented is the one that examines the interconnection between gold and oil prices and other econometric variables.

The results show that there were some causal connections among some of the electricity trading fundamentals which indicate that in order to facilitate the forecasting and minimize errors, we need to include them in the same model.

There are many different variables that can account for the changes in different financial markets. These changes may or may not always occur in financial markets. Nonetheless, this tutorial can demonstrate how we can use the relationship in certain financial markets to create models for forecasting.

Gold prices are closely related to other commodities. A hike in oil prices will have a positive impact on Gold prices and vice versa. Historically, we have seen that, when there is a hike in equity, Gold prices go down.

Varied sources of data had to be relied upon to collect the data related to Export and Import, WPI, and other Indices. All monthly data were collected from the said variables from EPWRF Data Set. These sources are vouched for

their authenticity in the data world, however, vary in the presentation of the data. Some use averages of the daily data readings to provide each monthly datum, while some just provide the data observed on the first or last day of the month. The data was not adjusted for such discrepancies, therefore is likely to have some bearing on the empirical results.

The WPI, Export, and Import Values were taken as proxies for Price Levels. To what extent do these proxies stand as a true representation of respective sectors is questionable itself and hence a limitation. Wholesale Price Indices only index a few types of consumer goods prices, may not be an accurate representation of the general price levels in the economy. This warrants consideration in future researches to use data sets that remain more representative of the testable parameters.

## References

1) Hiemstra, Craig, and Jonathan D. Jones. "Testing for linear and nonlinear Granger causality in the stock price- volume relation." The Journal of Finance 49.5 (1994): 1639-1664.

2) Johansen, S. (1995). Likelihood-based inference in cointegrated vector autoregressive models. Oxford University Press on Demand.

3) Lopez, J. Humberto. "The power of the ADF test." Economics Letters 57.1 (1997): 5-10.

4) Acharya, Rajesh H., and Anver C. Sadath. "Revisiting the relationship between oil price and macro economy: Evidence from India." ECONOMICS AND POLICY OF ENERGY AND THE ENVIRONMENT (2018).

## Code Listing

Check jupyter notebook, dataset and other details https://github.com/vntkumar8/
refactored-bassoon/.

```python
#!/usr/bin/env python
# coding: utf-8

# In[1]:


import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import warnings
import datetime


# In[2]:


df=pd.read_csv('Combined_data2.csv',parse_dates=['Time_Period'])


# In[3]:


df['Time_Period']=df['Time_Period'].apply(lambda x: datetime.datetime.
    strptime(x,'%b-%y'))


# In[4]:


df=df.set_index('Time_Period')


# In[5]:


df.head(5)


# In[169]:


from pandas.plotting import register_matplotlib_converters
register_matplotlib_converters()

# Plot
fig, axes = plt.subplots(nrows=3, ncols=3, dpi=120, figsize=(10,6))
for i, ax in enumerate(axes.flatten()):
    data = df[df.columns[i]]
    ax.plot(data, color='red', linewidth=1)
    # Decorations
    ax.set_title(df.columns[i])
    ax.xaxis.set_ticks_position('none')
    ax.yaxis.set_ticks_position('none')
    ax.spines['top'].set_alpha(0)
    ax.tick_params(labelsize=6)
    plt.tight_layout();

```

```
57
58  # In[170]:
59
60
61  from scipy import stats
62
63  # unpacks the test statistic and p-value from the normaltest of gold
64  stat, p = stats.normaltest(df.Oil_Price)
65
66  # prints both the test statistic and p-value
67  print('Statistics=%.3f, p=%.3f' % (stat, p))
68
69  # chooses an alpha of 0.05 and assigns to variable 'alpha'
70  alpha = 0.05
71
72  # creates an if statement for the hypothesis using the p-value test
73  if p > alpha:
74      print('Data looks Gaussian (fail to reject H0)')
75  else:
76      print('Data does not look Gaussian (reject H0)')
77
78
79  # In[172]:
80
81
82  print('Kurtosis of normal distribution: {}'.format(stats.kurtosis(df.
        Oil_Price)))
83  print('Skewness of normal distribution: {}'.format(stats.skew(df.Oil_Price)
        ))
84
85
86  # In[173]:
87
88
89  plt.figure(figsize=(14,6))
90  plt.subplot(1,2,1)
91  df['Oil_Price'].hist(bins=50)
92  plt.title('Oil')
93  plt.subplot(1,2,2)
94  stats.probplot(df['Oil_Price'], plot=plt);
95  df.Oil_Price.describe().T
96
97
98  # In[174]:
99
100
101 for name in X_test.columns:
102     plt.figure(figsize=(14,6))
103     plt.subplot(1,2,1)
104     df[name].hist(bins=50)
105     plt.title(name)
106     plt.subplot(1,2,2)
107     stats.probplot(df[name], plot=plt);
108     df.Oil_Price.describe().T
109
110
111 # In[179]:
112
113
114 for name in X_test.columns:
115     print(name)
116     # unpacks the test statistic and p-value from the normaltest of gold
```

```
117      stat, p = stats.normaltest(df[name])
118
119      # prints both the test statistic and p-value
120      print('Statistics=%.3f, p=%.3f' % (stat, p))
121
122      # chooses an alpha of 0.05 and assigns to variable 'alpha'
123      alpha = 0.05
124
125      # creates an if statement for the hypothesis using the p-value test
126      if p > alpha:
127          print('Data looks Gaussian (fail to reject H0)')
128      else:
129          print('Data does not look Gaussian (reject H0)')
130
131
132  # In[10]:
133
134
135  from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
136  from statsmodels.tsa.stattools import adfuller
137
138  fig, ax = plt.subplots(2, figsize=(10,6))
139  ax[0] = plot_acf(df[name], ax=ax[0])
140  ax[1] = plot_pacf(df[name], ax=ax[1])
141
142
143  # In[8]:
144
145
146  def adfuller_test(series, signif=0.05, name='', verbose=False):
147      """ Perform ADFuller to test for Stationarity of given series and print
          report"""
148      r = adfuller(series, autolag = 'BIC')
149      output = {'test_statistics':round(r[0], 4), 'pvalue':round(r[1], 4), '
          n_lags':round(r[2], 4), 'n_obs':r[3]}
150      p_value = output['pvalue']
151      def adjust(val, length=6): return str(val).ljust(length)
152
153      # Print Summary
154      print(f'    Augmented Dickey-Fuller Test on "{name}"', "\n   ", '-'*47)
155      print(f' Null Hypothesis: Data has unit root. Non-Stationary.')
156      print(f' Significance level    = {signif}')
157      print(f' Test Statistics       = {output["test_statistics"]}')
158      print(f' No. Lags Chosen       = {output["n_lags"]}')
159
160      for key, val in r[4].items():
161              print(f' Critical value {adjust(key)} = {round(val, 3)}')
162
163      if p_value <= signif:
164              print(f" => P-Value = {p_value}. Rejecting Null Hypothesis.")
165              print(f" => Series is Stationary.")
166
167      else:
168              print(f" => P-Value = {p_value}. Weak evidence to reject the Null
          Hypothesis.")
169              print(f" => Series in Non-Stationary.")
170
171
172  # In[11]:
173
174
175  for name, column in df.iteritems():
```

```
176            adfuller_test(column, name=column.name)
177            print('\n')
178
179
180 # In[12]:
181
182
183 nobs = 15
184
185 # assigning the last 15 observations to the variable 'X_test', while the
      rest will be assigned to the variable 'X_train'
186 X_train, X_test = df[0:-nobs], df[-nobs:]
187
188 #check size
189 print(X_train.shape)
190 print(X_test.shape)
191
192
193 # In[13]:
194
195
196 transform_data = X_train.diff().dropna()
197 transform_data.head()
198
199
200 # In[14]:
201
202
203 for name, column in transform_data.iteritems():
204            adfuller_test(column, name=column.name)
205            print('\n')
206
207
208 # In[15]:
209
210
211 fig, axes = plt.subplots(nrows=3, ncols=3, dpi=120, figsize=(12, 6))
212 for i, ax in enumerate(axes.flatten()):
213     data = transform_data[transform_data.columns[i]]
214     ax.plot(data, color='red', linewidth=1)
215     # Decorations
216     ax.set_title(transform_data.columns[i])
217     ax.xaxis.set_ticks_position('none')
218     ax.yaxis.set_ticks_position('none')
219     ax.spines['top'].set_alpha(0)
220     ax.tick_params(labelsize=6)
221
222 plt.tight_layout()
223
224
225 # In[16]:
226
227
228 # importing the granger causality test from statsmodels
229 from statsmodels.tsa.stattools import grangercausalitytests
230
231 # creating a variable named 'maxlag' and assigning the integer 5
232 maxlag = 5
233
234 # assigning the string 'ssr_chi2test' to the variable 'test'
235 test = 'ssr_chi2test'
236
```

```python
237  # creating a function named 'granger_causation_matrix' with the arguments '
         data', 'variables', 'test', and 'verbose'
238  def granger_causation_matrix(data, variables, test='ssr_chi2test', verbose=
         False):
239
240      # creating a dataframe with the same dimensions as number of variables
         entered, assigned to the variable 'X_train'
241      X_train = pd.DataFrame(np.zeros((len(variables), len(variables))),
         columns=variables, index=variables)
242
243      # loops through the columns and the indexes
244      for c in X_train.columns:
245          for r in X_train.index:
246
247              # conducts a granger causality test on a variable row and
         column using the 'maxlag' variable; assigns to
248              # variable test_result
249              test_result = grangercausalitytests(data[[r, c]], maxlag=maxlag
         , verbose=False)
250
251              # locates the test result in the tuple 'test_result' and rounds
          the number by 4 digits; assigns to 'p_values'
252              p_values = [round(test_result[i+1][0][test][1], 4) for i in
         range(maxlag)]
253
254              # if the variable is 'verbose', print the y and x variables and
          the associated p-value
255              if verbose:
256                  print(f'Y = {r}, X = {c}, P Values = {p_values}')
257
258              # find the smallest p-value and assign to variable 'min_p_value
         '
259              min_p_value = np.min(p_values)
260
261              # find the smallest p-value and assign to its respective row
         and column
262              X_train.loc[r, c] = min_p_value
263
264      # remane the row and column names based on the relationship
265      X_train.columns = [var + '_x' for var in variables]
266      X_train.index = [var + '_y' for var in variables]
267      return X_train
268
269
270  # In[17]:
271
272
273  granger_causation_matrix(X_train, variables = X_train.columns)
274
275
276  # In[25]:
277
278
279  from statsmodels.tsa.vector_ar.vecm import coint_johansen
280
281  def cointegration_test(transfrom_data, alpha=0.05):
282      """Perform Johanson's Cointegration Test and Report Summary"""
283      for i in range(0,8):
284          print("for i="+str(i)+'\n')
285          out = coint_johansen(transform_data, -1, i)
286          d = {'0.90':0, '0.95':1, '0.99':2}
287          traces = out.lr1
```

```
288          cvts = out.cvt[:, d[str(1-alpha)]]
289          def adjust(val, length=6):
290              return str(val).ljust(length)
291
292          # Summary
293          print('Name   :: Test Stat > C(95%)       => Signif  \n', '--'*20)
294          for col, trace, cvt in zip(transform_data.columns, traces, cvts):
295              print(adjust(col), '::', adjust(round(trace, 2), 9), ">",
     adjust(cvt, 8), ' =>  ', trace >cvt)
296 #      print(traces)
297 #      print(cvts)
298
299 cointegration_test(X_train)
300
301
302 # In[191]:
303
304
305 import statsmodels.tsa.api as smt
306 from statsmodels.tsa.api import VAR
307 mod = smt.VAR(transform_data)
308 res = mod.fit(maxlags= 5, ic = 'aic')
309 print(res.summary())
310
311
312 # In[192]:
313
314
315 pred = res.forecast(transform_data.values[-9:], 15)
316 pred_df = pd.DataFrame(pred, index=df.index[-15:], columns = df.columns)
317 pred_df
318
319
320 # In[193]:
321
322
323 # reversing the difference and assigning to variable 'pred_inverse'
324 pred_inverse = pred_df.cumsum()
325
326 # inverse the difference values and assigning to variable 'f'
327 f = pred_inverse + X_test
328 print(f)
329
330
331 # In[194]:
332
333
334 plt.figure(figsize=(12,5))
335 plt.xlabel('Date')
336
337 ax1 = X_test.Gold_Price.plot(color='blue', grid=True, label='Actual Gold
     Price')
338 ax2 = f.Gold_Price.plot(color='red', grid=True, secondary_y=True, label='
     Predicted Gold Price')
339
340 ax1.legend(loc=1)
341 ax2.legend(loc=2)
342 plt.title('Prediction vs Actual Gold Price')
343 plt.show()
344
345
346 # In[195]:
```

```
347
348
349 plt.figure(figsize=(12,5))
350 plt.xlabel('Date')
351
352 for name in X_test.columns:
353     label=str(name)
354     ax1 = X_test[name].plot(color='blue', grid=True, label='Actual')
355     ax2 = f[name].plot(color='red', grid=True, secondary_y=True, label='
        Predicted')
356
357     ax1.legend(loc=1)
358     ax2.legend(loc=2)
359     plt.title(label)
360     plt.show()
361
362
363 # In[115]:
364
365
366 plt.figure(figsize=(12,5))
367 plt.xlabel('Date')
368
369 ax1 = X_test.Oil_Price.plot(color='blue', grid=True, label='Actual Oil
        Price')
370 ax2 = f.Oil_Price.plot(color='red', grid=True, secondary_y=True, label='
        Predicted Oil Price')
371
372 ax1.legend(loc=1)
373 ax2.legend(loc=2)
374 plt.title('Prediction vs Actual Oil Price')
375 plt.show()
376
377
378 # In[205]:
379
380
381 from sklearn.metrics import mean_absolute_error
382 from sklearn.metrics import mean_squared_error
383 from math import sqrt
384
385 forecast_errors = [X_test.Oil_Price[i]-f.Oil_Price[i] for i in range(len(
        X_test.Oil_Price))]
386 bias = sum(forecast_errors) * 1.0/len(X_test.Oil_Price)
387 print('Bias: %f' % bias)
388
389 mae = mean_absolute_error(X_test.Oil_Price, f.Oil_Price)
390 print('MAE: %f' % mae)
391
392 mse = mean_squared_error(X_test.Oil_Price, f.Oil_Price)
393 print('MSE: %f' % mse)
394
395 rmse = sqrt(mse)
396 print('RMSE: %f' % rmse)
397
398
399 # In[207]:
400
401
402 for name in X_test.columns:
403     print(name)
```

```
404    forecast_errors = [X_test[name][i]-f[name][i] for i in range(len(X_test
       [name]))]
405    bias = sum(forecast_errors) * 1.0/len(X_test[name])
406    print('Bias: %f' % bias)
407
408    mae = mean_absolute_error(X_test[name], f[name])
409    print('MAE: %f' % mae)
410
411    mse = mean_squared_error(X_test[name], f[name])
412    print('MSE: %f' % mse)
413
414    rmse = sqrt(mse)
415    print('RMSE: %f' % rmse)
416    print('-'*10)
```