

DSC 275/475: Time Series Analysis and Forecasting (Fall 2022)

Project-1

(Total points: 60 for undergraduate students; 70 for Graduate students and including extra credit for undergraduate students)

Aradhya Mathur

Overview This project is designed to provide you hands-on experience working on an end-to-end time series analysis and forecasting solution using AR/ARMA/ARIMA/SARIMA modeling. You are welcome to use any external libraries/packages for this project. A few recommendations are provided for each problem below. This is a guided exercise in that you are expected to answer each of the questions below. For some of the questions you will appreciate there is no single correct answer. In such cases, you have flexibility to decide the approach. Any conclusions you make or decisions you take must be stated and accompanied by a reasonable justification. Your submission should be a PDF document with responses (including figures/plots) to each question along with the code either included inline [e.g. Notebook] or as a separate file.

In [12]:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from statsmodels.tsa.stattools import adfull, diff
from statsmodels.tsa.arima.model import ARIMA
import statsmodels.graphics.tsaplots as tsa
import warnings
warnings.filterwarnings("ignore")
from statsmodels.tsa.stattools import acf, pacf
from sklearn.metrics import mean_absolute_error
from sklearn.metrics import mean_squared_error
```

Question 1

Problem 1 (60 pts; Required for all students) The data for this project (Problem1_DataSet.csv) represents 7 years of monthly data on airline miles flown in the United Kingdom. You are tasked with the goal of developing a forecasting model that can accurately predict the trend for future years. To achieve the final goal, answer each of the questions below

1. Create a time series of the plot of the data provided. (5 pts)

```
# Reading dataset
dfpl = pd.read_csv('Problem1_DataSet.csv', header = 0)
dfpl.head()
```

Out[31]:

```
Month  Miles, in Millions
0  Jan-1964      7.269
1  Feb-1964      6.775
2  Mar-1964      7.819
3  Apr-1964      8.371
4  May-1964      9.069
```

In [32]:

```
#Changing Month format
dfpl['Month'] = pd.to_datetime(dfpl['Month'])
dfpl.head()
```

Out[32]:

```
Month  Miles, in Millions
0  1964-01-01      7.269
1  1964-02-01      6.775
2  1964-03-01      7.819
3  1964-04-01      8.371
4  1964-05-01      9.069
```

In [33]:

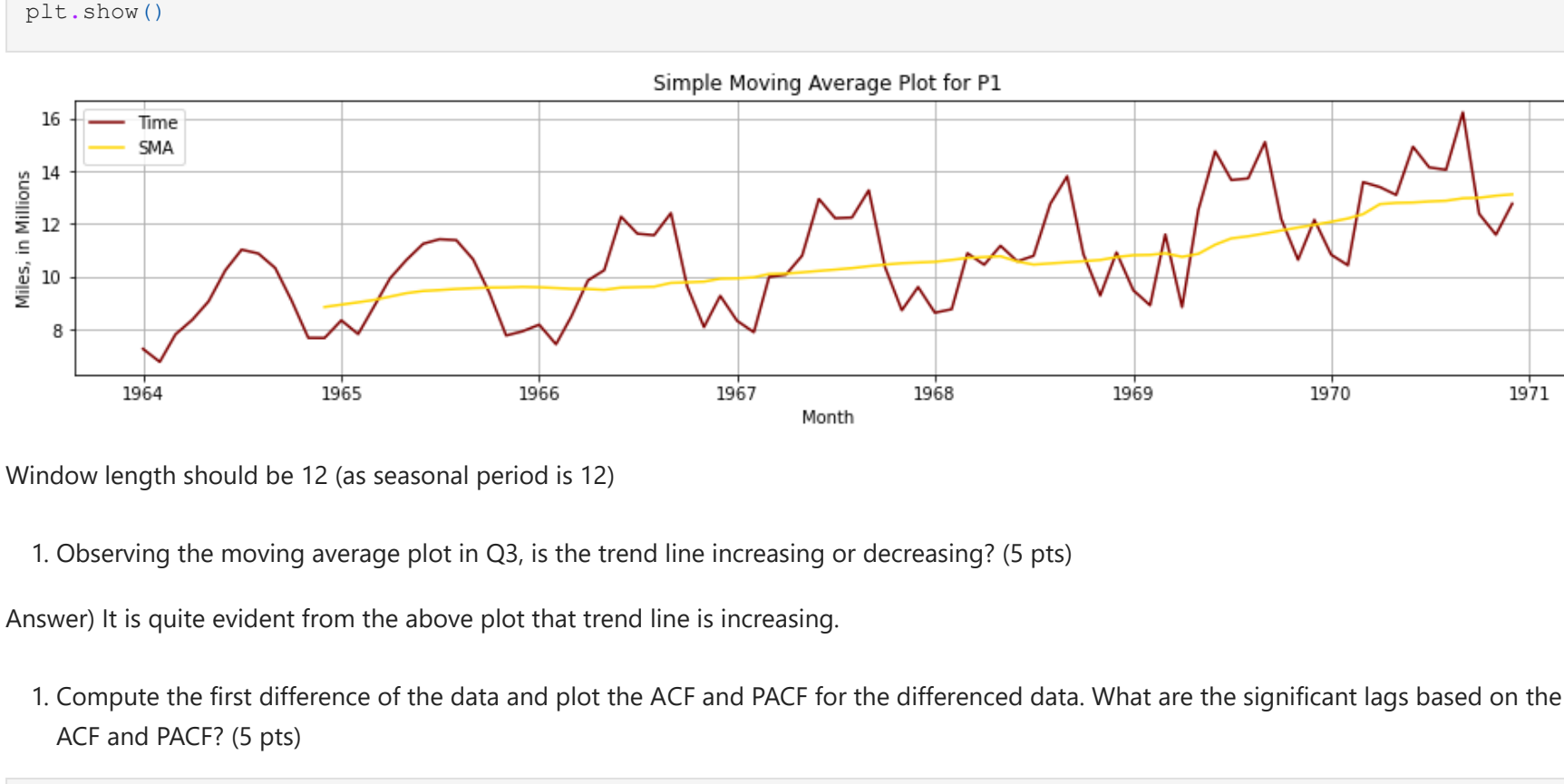
```
#Renaming Column
dfpl.columns = ['Month', 'Miles']
dfpl.head()
```

Out[33]:

```
Month  Miles
0  1964-01-01      7.269
1  1964-02-01      6.775
2  1964-03-01      7.819
3  1964-04-01      8.371
4  1964-05-01      9.069
```

In [34]:

```
#Time Series Plot
plt.figure(figsize=(16, 3))
plt.plot(dfpl['Month'], dfpl['Miles'], label = 'Time Series Plot', color = 'maroon')
plt.xlabel('Month')
plt.ylabel('Miles, in Millions')
plt.legend()
plt.title('Time Series Plot - Miles vs Miles')
plt.show()
```



2. Plot the autocorrelation function (ACF). From the ACF, what is the seasonal period? (5 pts)

```
# ACF Plot
plt.figure(figsize=(16, 3))
tsa.plot_acf(dfpl['Miles'], lags = 50, color = 'maroon')
plt.xlabel('Lag')
plt.ylabel('Autocorrelation')
plt.title('ACF Plot')
plt.show()
```

Figure size 1152x216 with 0 Axes



Seasonal Period is 12 lags (0-11 first period)

1. Compute a moving average for the data to determine the trend in the data and overlay on the original time-series plot. What is a suitable choice for the moving average window length? (5 pts)

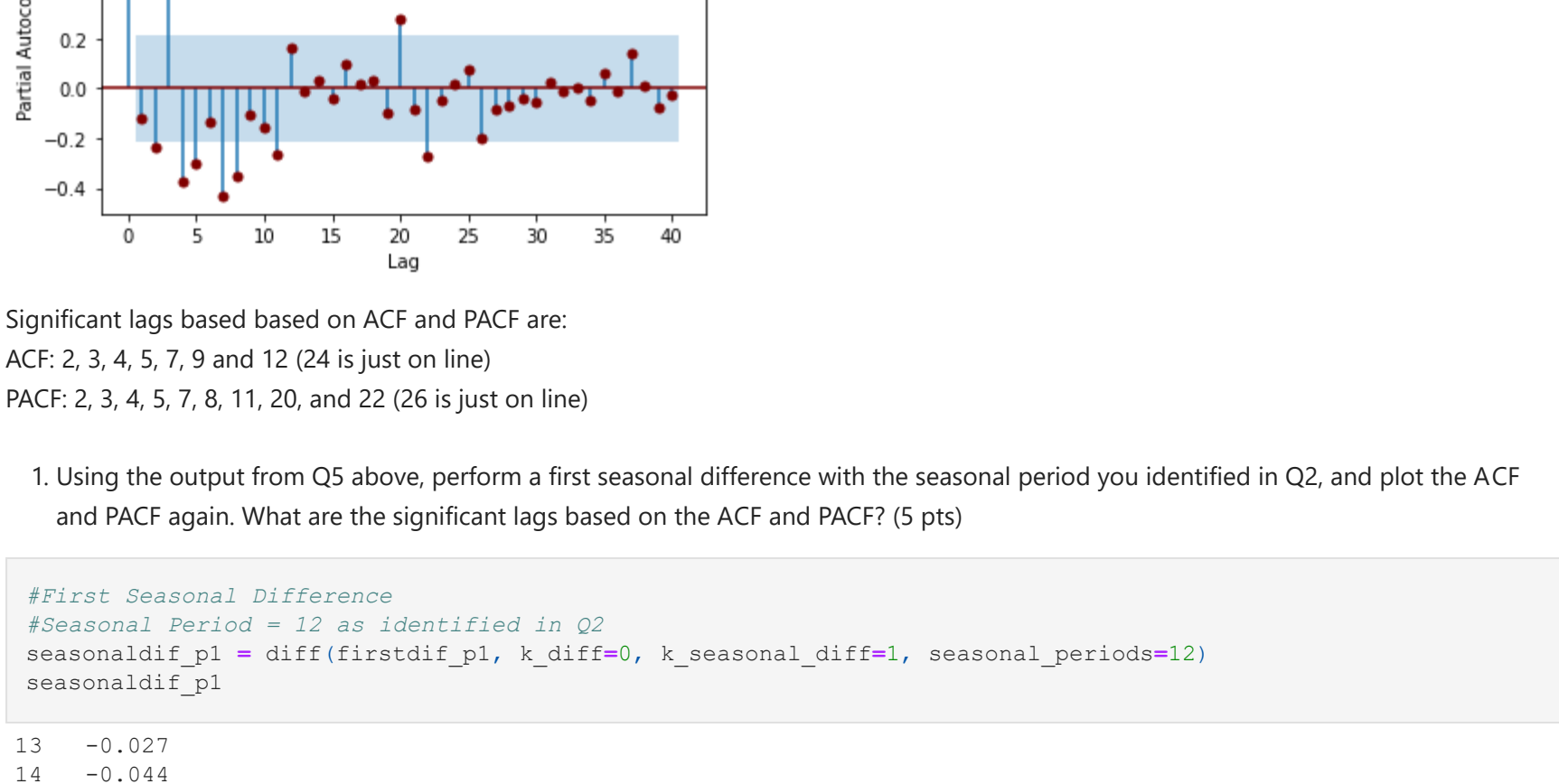
```
#Simple Moving Average
dfpl['Sma_p1'] = dfpl['Miles'].rolling(12, min_periods=12).mean()
dfpl.head()
```

Out[36]:

```
Month  Miles  Sma_p1
0  1964-01-01      NaN
1  1964-02-01      NaN
2  1964-03-01      NaN
3  1964-04-01      NaN
4  1964-05-01      NaN
```

In [37]:

```
#Overlaying Time Series Plot
plt.figure(figsize=(16, 3))
monthpl = dfpl['Month']
milespl = dfpl['Miles']
plt.plot(monthpl, milespl, label = 'Time', color = 'maroon')
plt.plot(monthpl, dfpl['Sma_p1'], label = 'SMA', color = 'gold')
plt.xlabel('Month')
plt.ylabel('Miles, in Millions')
plt.legend()
plt.show()
```



Window length should be 12 (as seasonal period is 12)

1. Observing the moving average plot in Q3, is the trend line increasing or decreasing? (5 pts)

Answer) It is quite evident from the above plot that trend line is increasing.

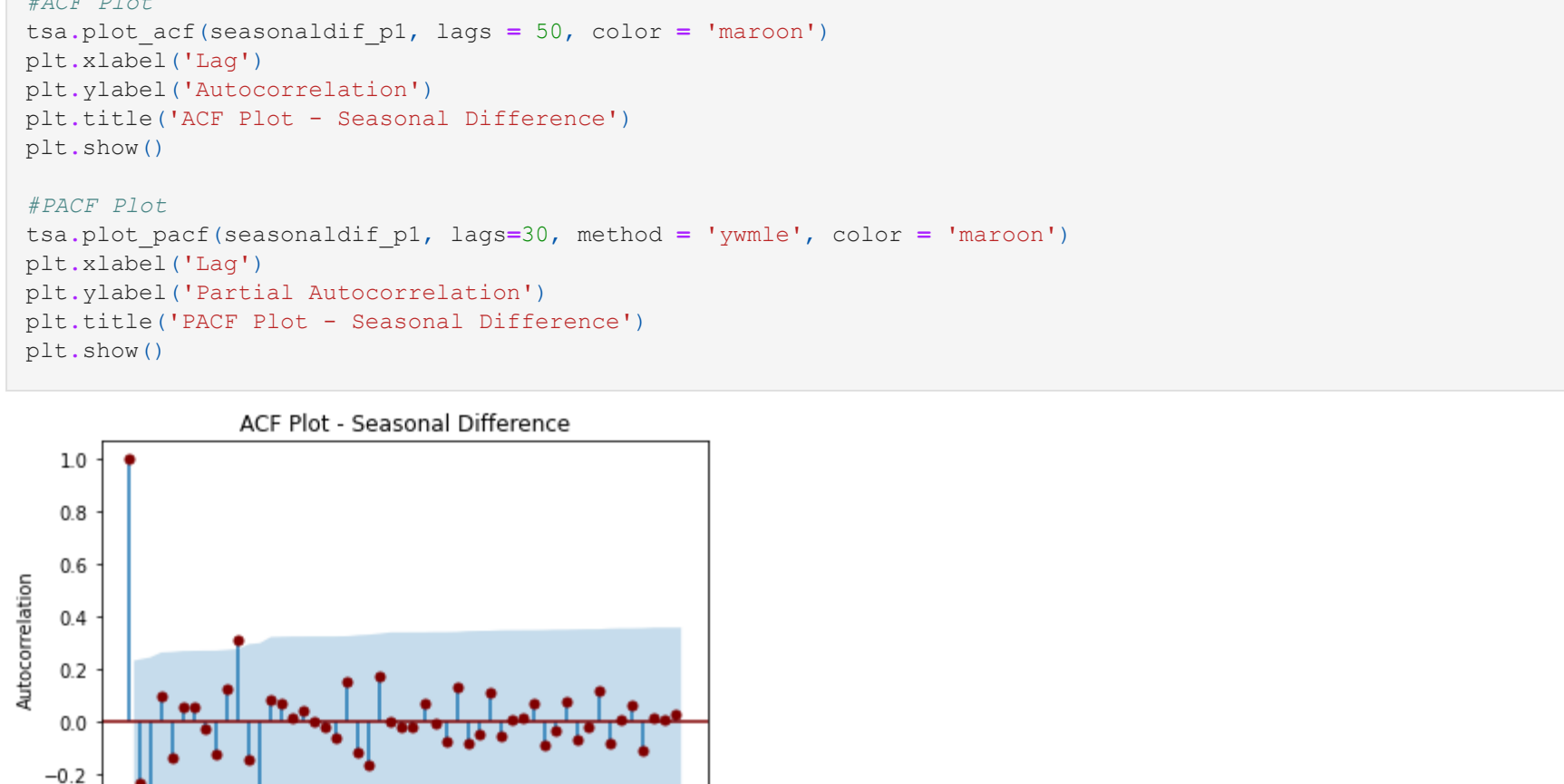
1. Compute the first difference of the data and plot the ACF and PACF for the differenced data. What are the significant lags based on the ACF and PACF? (5 pts)

```
#First difference
firstdif_pl = diff(milespl)
print(firstdif_pl)
```

```
1    -0.494
2     1.044
3     0.552
4     0.698
5     1.179
...
79    -0.090
80     2.177
81    -3.845
82    -0.798
83     1.178
Name: Miles, Length: 83, dtype: float64
```

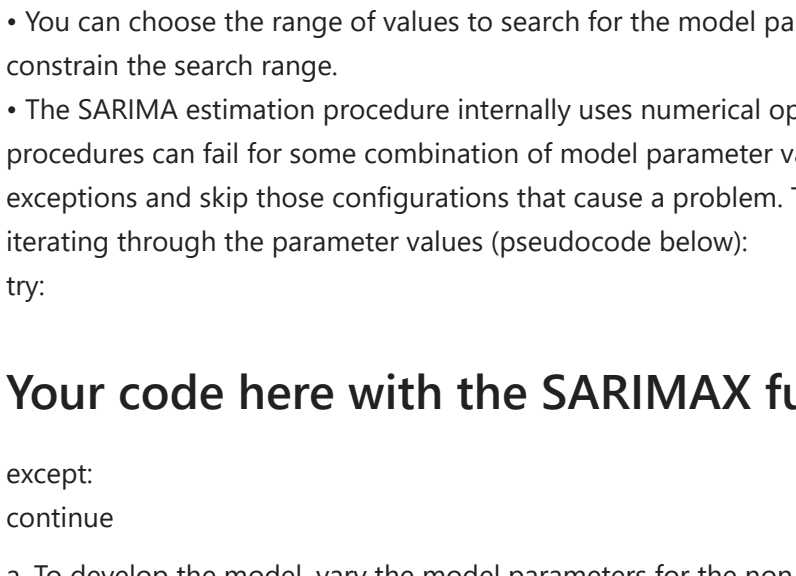
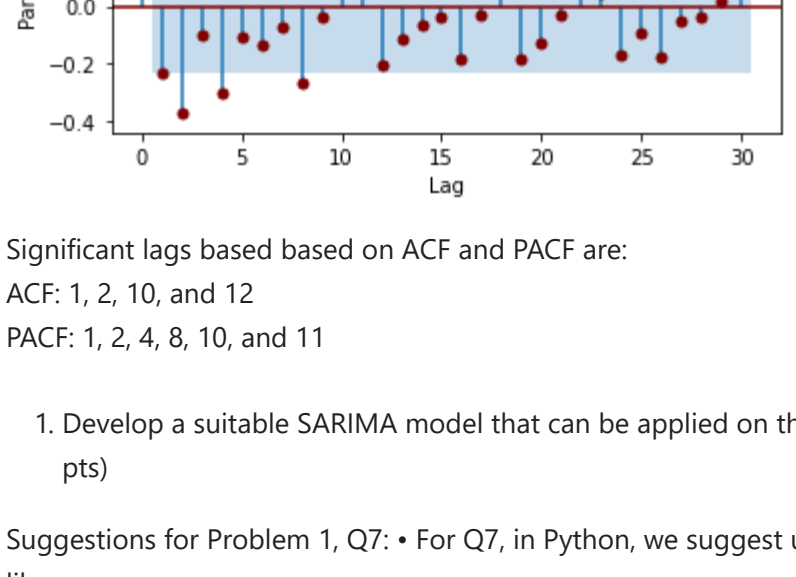
In [39]:

```
#Plotting first difference
plt.figure(figsize=(16, 3))
firstdif_pl.plot(label = 'First Difference', color = 'maroon')
plt.grid()
plt.xlabel('Month')
plt.ylabel('First Difference')
plt.legend()
plt.title('Plot - First Difference')
plt.show()
```



In [40]:

```
#ACF of first difference
tsa.plot_acf(firstdif_pl, lags = 50, color = 'maroon')
plt.figure(figsize=(16, 3))
plt.xlabel('Lag')
plt.ylabel('Autocorrelation')
plt.title('ACF Plot- First Difference')
plt.legend()
#PACF of first difference
tsa.plot_pacf(firstdif_pl, lags=40, method = 'ywmls', color = 'maroon') # Not possible for lag = 50
plt.xlabel('Lag')
plt.ylabel('Partial Autocorrelation')
plt.title('PACF Plot- First Difference')
plt.show()
```



Significant lags based on ACF and PACF are:

ACF: 2, 3, 4, 5, 7, 9 and 12 (24 is just on line)

PACF: 2, 3, 4, 5, 7, 8, 11, 20, and 22 (26 is just on line)

1. Using the output from Q5 above, perform a first seasonal difference with the seasonal period you identified in Q2, and plot the ACF and PACF again. What are the significant lags based on the ACF and PACF? (5 pts)

```
#First Seasonal Difference
#We choose seasonal period = 12 as identified in Q2
seasonaldif_pl = diff(firstdif_pl, k_diff=0, k_seasonal_diff=1, seasonal_periods=12)
seasonaldif_pl
```

```
13    -0.027
14    -0.044
15    -0.567
16    -0.008
17    -0.564
...
79    -0.154
80    -0.798
81    -0.920
82    -0.745
83    -0.338
Name: Miles, Length: 71, dtype: float64
```

In [42]:

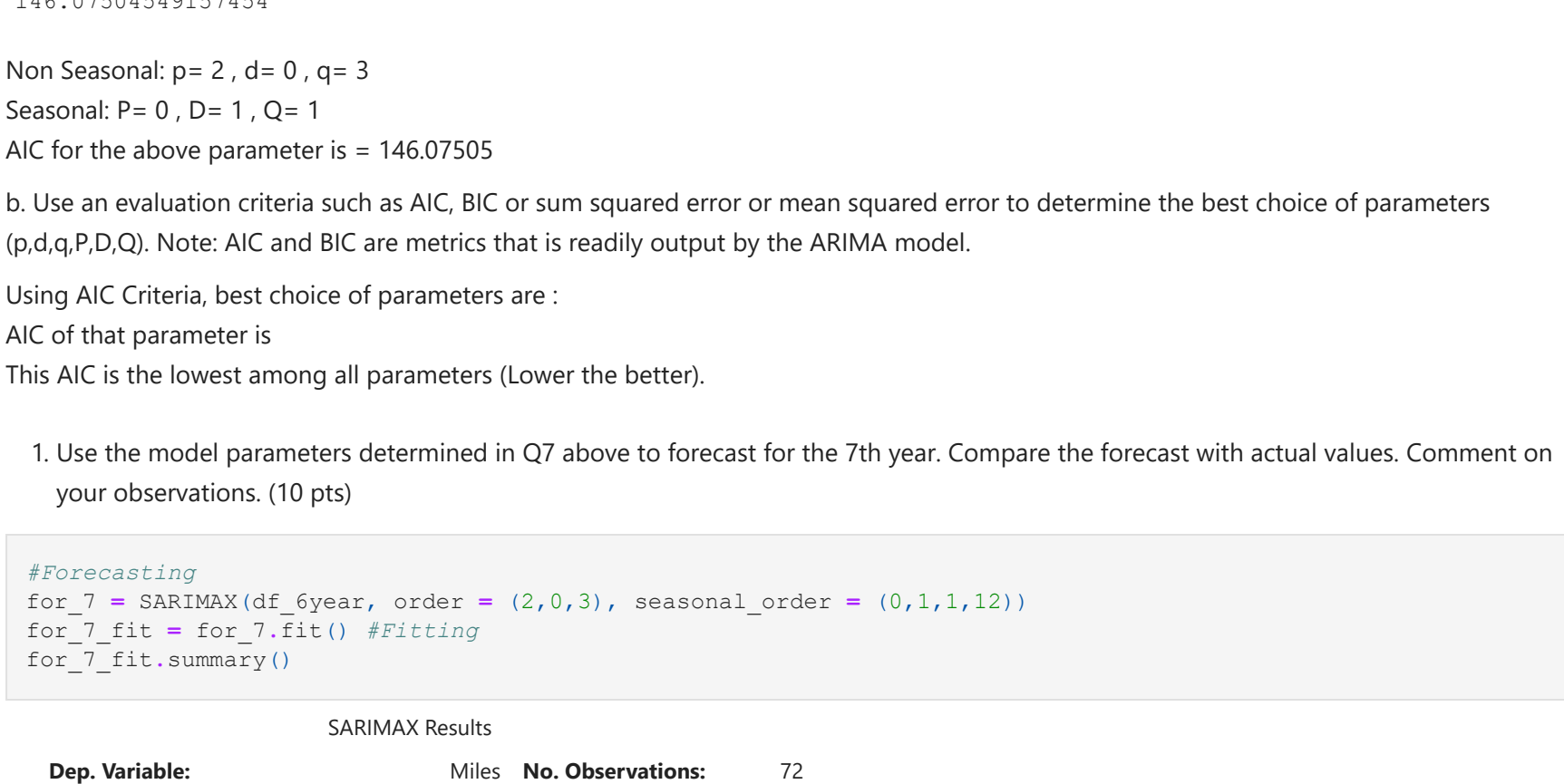
```
#Checking if on original data, we apply k_diff = 1 which is what we did in question 5
seasonaldif_p2 = diff(milespl, k_diff=1, k_seasonal_diff=1, seasonal_periods=12)
seasonaldif_p2
#Output is same
```

Out[42]:

```
13    -0.027
14    -0.044
15    -0.567
16    -0.008
17    -0.564
...
79    -0.154
80    -0.798
81    -0.920
82    -0.745
83    -0.338
Name: Miles, Length: 71, dtype: float64
```

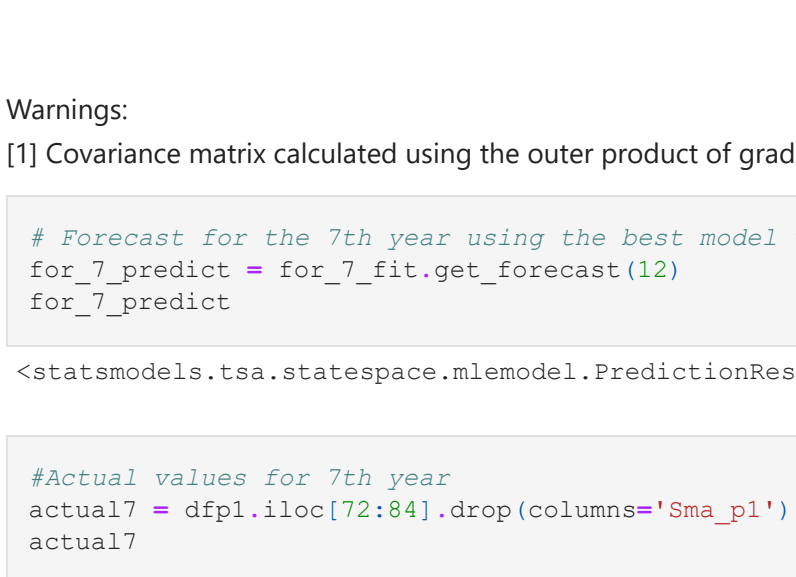
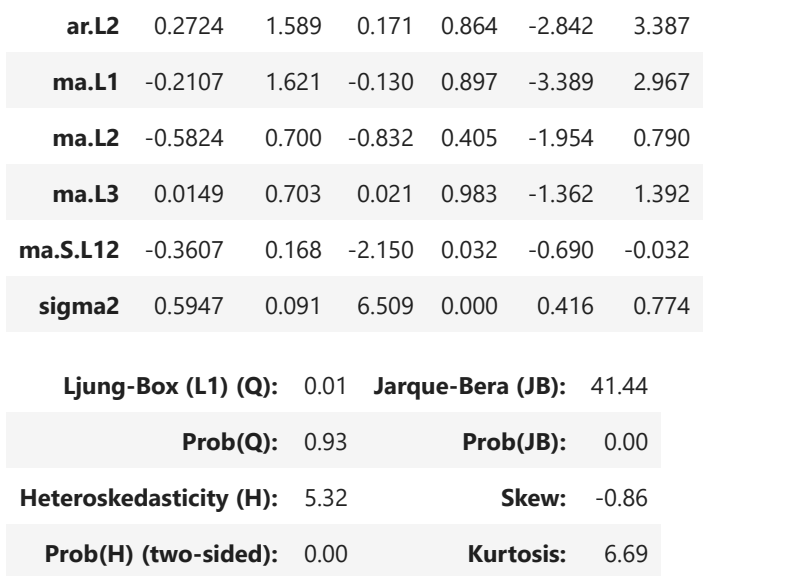
In [43]:

```
#Plotting first difference
plt.figure(figsize=(16, 3))
seasonaldif_pl.plot(label = 'First Seasonal Difference', color = 'maroon')
plt.grid()
plt.xlabel('Month')
plt.ylabel('First Seasonal Difference')
plt.legend()
plt.title('Plot - First Seasonal Difference')
plt.show()
```



In [44]:

```
#ACF Plot
tsa.plot_acf(seasonaldif_pl, lags = 50, color = 'maroon')
plt.xlabel('Lag')
plt.ylabel('Autocorrelation')
plt.title('ACF Plot - Seasonal Difference')
plt.show()
#PACF Plot
tsa.plot_pacf(seasonaldif_pl, lags=30, method = 'ywmls', color = 'maroon')
plt.xlabel('Lag')
plt.ylabel('Partial Autocorrelation')
plt.title('PACF Plot - Seasonal Difference')
plt.show()
```



Significant lags based on ACF and PACF are:

ACF: 1, 2, 10, and 12

PACF: 1, 2, 4, 8, 10, and 11

1. Develop a suitable SARIMA model that can be applied on the time series. Use the first 6 years of data only to develop the model. (20 pts)

Suggestions for Problem 1, Q7: • For Q7, in Python, we suggest using the package/function SARIMAX in the 'statsmodels.tsa.statespace' library

• You can choose the range of values to search for the model parameters. We suggest varying p, q and P, Q each over the range 0 to 3 to constrain the search range.

• The SARIMA estimation procedure internally uses numerical optimization procedures to find a set of coefficients for the model. These procedures can fail for some combination of model parameter values which in turn can throw Python errors. We must catch these exceptions and skip those combinations that cause a problem. To solve this problem, include 'try/except' blocks in your code when iterating through the parameter values (pseudocode below):

try:

Your code here with the SARIMAX function

except:

continue

a. To develop the model, vary the model parameters for the non-seasonal (p,q,d) and seasonal components (P,D,Q) and calculate the output for each combination of parameters.

```
from statsmodels.tsa.statespace.sarimax import SARIMAX
```

In [46]:

```
#Using first 6 years of data
df_6year = dfpl['Miles'].iloc[0:72] # 6 years = 72 Months
df_6year
```

Out[46]:

```
0    7.269
1    6.775
2    7.819
3    8.371
4    9.069
...
67   13.731
68   15.110
69   12.185
70   10.645
71   12.161
Name: Miles, Length: 72, dtype: float64
```

In [47]:

```
from itertools import product
```

In [97]:

```
#Building Model and DataFrame
df_sarima_model = pd.DataFrame()
min_aic = []
for [p,d,q,P,D,Q] in product([0,1,2,3],[0,1],[0,1,2,3],[0,1],[0,1,2,3],[0,1],[0,1,2,3]):
    try:
        m = SARIMAX(df_6year, trend='c', order = (p,d,q), seasonal_order = (P,D,Q,12)) # Sarima Model
        AIC = model.aic # AIC calculation
        min_aic.append(AIC) # Finding minimum AIC
        min_aic.sort()
        df_sarima_model = df_sarima_model.append({'p': p, 'd': d, 'q': q, 'P': P, 'D': D, 'Q': Q, 'AIC': AIC})
    except:
        continue
```

```
#Finding Minimum AIC
print("Minimum AIC is =", min_aic[0])
```

Minimum AIC is = 22.0

In [103]:

```
df_sarima_model.head()
```

Out[103]:

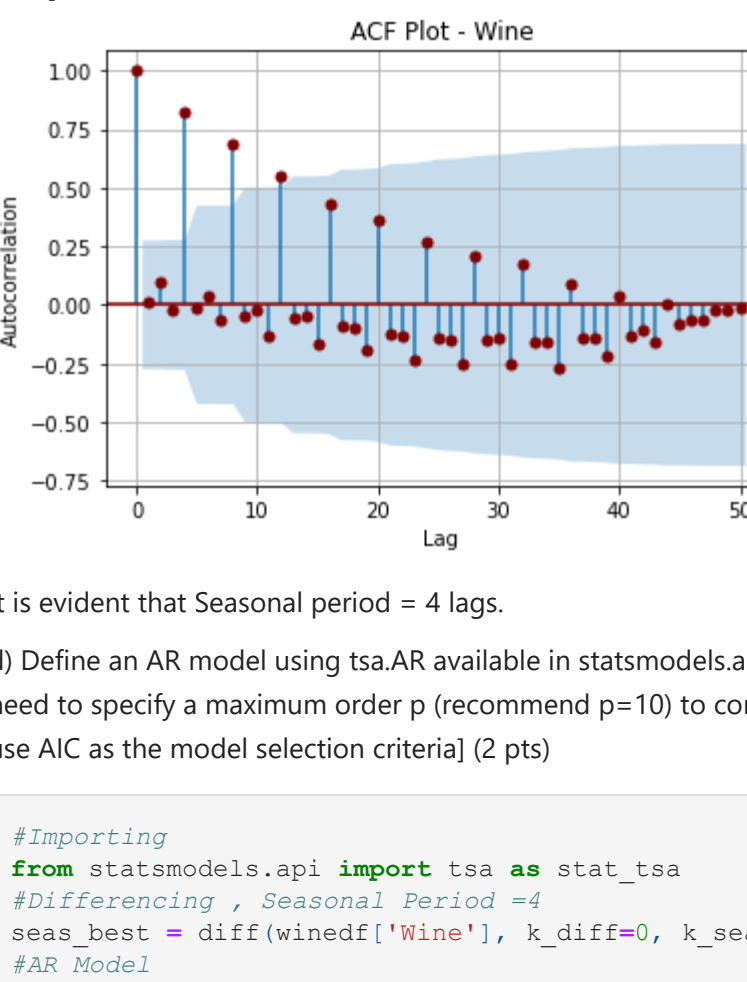
```
p  d  q  P  D  Q  AIC
0  0  0  0  0  0  0  300.658661
1  0  0  0  0  0  0  10  272.463181
2  0  0  0  0  0  0  20  256.994615
3  0  0  0  0  0  0  30  248.434514
4  0  0  0  0  0  0  40  164.870606
```

In [104]:

```
df_sarima_model.sort_values(by=['AIC']).head()
```

```
p  d  q  P  D  Q  AIC
847  0  0  0  2  0  10  30  22.000000
848  0  0  0  3  0  10  10  22.000000
849  0  0  0  4  0  10  10  22.000000
850  0  0  0  5  0  10  10  22.000000
851  0  0  0  6  0  10  10  22.000000
852  0  0  0  7  0  10  10  22.000000
853  0  0  0  8  0  10  10  22.000000
854  0  0  0  9  0  10  10  22.000000
855  0  0  0  10  0  10  10  22.000000
856  0  0  0  11  0  10  10  22.000000
857  0  0  0  12  0  10  10  22.000000
858  0  0  0  13  0  10  10  22.000000
859  0  0  0  14  0  10  10  22.000000
860  0  0  0  15  0  10  10  22.000000
861  0  0  0  16  0  10  10  22.000000
862  0  0  0  17  0  10  10  22.000000
863  0  0  0  18  0  10  10  22.000000
864  0  0  0  19  0  10  10  22.000000
865  0  0  0  20  0  10  10  22.000000
866  0  0  0  21  0  10  10  22.000000
867  0  0  0  22  0  10  10  22.000000
868  0  0  0  23  0  10  10  22.000000
869  0  0  0  24  0  10  10  22.000000
870  0  0  0  25  0  10  10  22.000000
871  0  0  0  26  0  10  10  22.000000
872  0  0  0  27  0  10  10  22.000000
873  0  0  0  28  0  10  10  22.000000
874  0  0  0  29  0  10  10  22.000000
875  0  0  0  30  0  10  10  22.000000
876  0  0  0  31  0  10  10  22.000000
877  0  0  0  32  0  10  10  22.000000
878  0  0  0  33  0  10  10  22.000000
879  0  0  0  34  0  10  10  22.000000
880  0  0  0  35  0  10  10  22.000000
881  0  0  0  36  0  10  10  22.000000
882  0  0  0  37  0  10  10  22.000000
883  0  0  0  38  0  10  10  22.000000
884  0  0  0  39  0  10  10  22.000000
885  0  0  0  40  0  10  10  22.000000
886  0  0  0  41  0  10  10  22.000000
887  0  0  0  42  0  10  10  22.000000
888  0  0  0  43  0  10  10  22.000000
889  0  0  0  44  0  10  10  22.000000
890  0  0  0  45  0  10  10  22.000000
891  0  0  0  46  0  10  10  22.000000
892  0  0  0  47  0  10  10  22.000000
893  0  0  0  48  0  10  10  22.000000
894  0  0  0  49  0  10  10  22.000000
895  0  0  0  50  0  10  10  22.000000
896  0  0  0  51  0  10  10  22.000000
897  0  0  0  52  0  10  10  22.000000
898  0  0  0  53  0  10  10  22.000000
899  0  0  0  54  0  10  10  22.000000
900  0  0  0  55  0  10  10  22.000000
901  0  0  0  56  0  10  10  22.000000
902  0  0  0  57  0  10  10  22.000000
903  0  0  0  58  0  10  10  22.000000
904  0  0  0  59  0  10  10  22.000000
905  0  0  0  60  0  10  10  22.000000
906  0  0  0  61  0  10  10  22.000000
907  0  0  0  62  0  10  10  22.000000
908  0  0  0  63  0  10  10  22.000000
909  0  0  0  64  0  10  10  22.000000
910  0  0  0  65  0  10  10  22.000000
911  0  0  0  66  0  10  10  22.000000
912  0  0  0  67  0  10  10  22.000000
913  0  0  0  68  0  10  10  22.000000
914  0  0  0  69  0  10  10  22.000000
915  0  0  0  70  0  10  10  22.000000
916  0  0  0  71  0  10  10  22.000000
917  0  0  0  72  0  10  10  22.000000
918  0  0  0  73  0  10  10  22.000000
919  0  0  0  74  0  10  10  22.000000
920  0  0  0  75  0  10  10  22.000000
921  0  0  0  76  0  10  10  22.000000
922  0  0  0  77  0  10  10  22.000000
923  0  0  0  78  0  10  10  22.000000
924  0  0  0  79  0  10  10  22.000000
925  0  0  0  80  0  10  10  22.000000
926  0  0  0  81  0  10  10  22.000000
927  0  0  0  82  0  10  10  22.000000
928  0  0  0  83  0  10  10  22.000000
929  0  0  0  84  0  10  10  22.000000
930  0  0  0  85  0  10  10  22.000000
931  0  0  0  86  0  10  10  22.000000
932  0  0  0  87  0  10  10  22.000000
933  0  0  0  88  0  10  10  22.000000
934  0  0  0  89  0  10  10  22.000000
935  0  0  0  90  0  10  10  22.000000
936  0  0  0  91  0  10  10  22.000000
937  0  0  0  92  0  10  10  22.000000
938  0  0  0  93  0  10  10  22.000000
939  0  0  0  94  0  10  10  22.000000
940  0  0  0  95  0  10  10  22.000000
941  0  0  0  96  0  10  10  22.000000
942  0  0  0  97  0  10  10  22.000000
943  0  0  0  98  0  10  10  22.000000
944  0  0  0  99  0  10  10  22.000000
945  0  0  0  100  0  10  10  22.000000
946  0  0  0  101  0  10  10  22.000000
947  0  0  0  102  0  10  10  22.000000
948  0  0  0  103  0  10  10  22.000000
949  0  0  0  104  0  10  10  22.000000
950  0  0  0  105  0  10  10  22.000000
951  0  0  0  106  0  10  10  22.000000
952  0  0  0  107  0  10  10  22.000000
953  0  0  0  108  0  10  10  22.000000
954  0  0  0  109  0  10  10  22.000000
955  0  0  0  110  0  10  10  22.000000
956  0  0  0  111  0  10  10  22.000000
957  0  0  0  112  0  10  10  22.000000
958  0  0  0  113  0  10  10  22.000000
959  0  0  0  114  0  10  10  22.000000
960  0  0  0  115  0  10  10  22.000000
961  0  0  0  116  0  10  10  22.000000
962  0  0  0  117  0  10  10  22.000000
963  0  0  0  118  0  10  10  22.000000
964  0  0  0  119  0  10  10  22.000000
965  0  0  0  120  0  10  10  22.000000
966  0  0  0  121  0  10  10  22.000000
967  0  0  0  122  0  10  10  22.000000
968  0  0  0  123  0  10  10  22.000000
969  0  0  0  124  0  10  10  22.000000
970  0  0  0  125  0  10  10  22.000000
971  0  0  0  126  0  10  10  22.000000
972  0  0  0  127  0  10  10  22.000000
973  0  0  0  128  0  10  10  22.000000
974  0  0  0  129  0  10  10  22.000000
975  0  0  0  130  0  10  10  22.000000
976  0  0  0  131  0  10  10  22.000000
977  0  0  0 
```


<Figure slice 1152x2 with 0 Axes>



It is evident that Seasonal period = 4 lags.

d) Define an AR model using tsa.AR available in statsmodels.api. Determine the optimal order using the "select_order" function. You will need to specify a maximum order p (recommend p=10) to consider and a criterion for deciding which model order is "best". (e.g. You can use AIC as the model selection criteria) [2 pts]

```
In [121]: #Importing
from statsmodels.api import tsa as stat_tsa
#Differencing , Seasonal Period =4
seas_best = diff(winedf['Wine'], k_diff=0, k_seasonal_diff=1, seasonal_periods=4)
#AR Model
ar_wine = stat_tsa.AR(seas_best)
#Finding optimal order (p=10) using select_order
optimal = ar_wine.select_order(maxlag=10, ic='aic')
optimal
```

Out[121]: 5

Optimal order is 5

e) Now, evaluate an AR(p) model for the time-series generated after seasonal differencing (using the best lag you found in part b above) (4 pts)

i. use the fit method specifying the optimal lag found above

```
In [122]: #Fit method for optimal lag = 5
ar_model = stat_tsa.AR(seas_best)
ar_model_fit = ar_model.fit(maxlag = 5)
ar_model_fit.summary()
```

```

Wine_pred[3:10, :].head()

9      0.057
10     0.231
11     0.163
12     0.035
13     0.115
dtype: float64

# Using predict method to generate values
predict_ar = ar_model_fit.predict(start = 5)
predict_ar[0:42] = predict_ar[0:42] # length = 42
predict_ar.head()

9      0.146767
10     -0.022745
11     0.217034
12     -0.137246
13     0.136466
dtype: float64

# plot the predicted results and the corresponding seasonally

#Plot the Predicted Results and the Corresponding
plt.figure(figsize=(15, 3))
# Predict values
plt.grid()
plt.plot(predict_ar, label = 'Predicted Time Se
# Seasonally Differenced Time Series
plt.plot(seas_best, label = 'Seasonally Differ
# Predicted results and the corresponding
plt.xlabel('Time (Quarter)')
plt.ylabel('Seasonal Difference')
plt.legend()
plt.show()

```

ii. use the predict method to generate values starting at the optimal lag

```
In [123]: seas_best[5:].count()
```

Out[123]: 42

```
In [124]: seas_best[5:].head()
```

Out[124]: 9 0.057
10 0.231
11 0.163
12 0.035
13 0.115
Name: Wine, dtype: float64

```
In [125]: # Using predict method to generate values starting at the optimal lag = 5
predict_ar = ar_model_fit.predict(start = 5)
predict_ar_p = predict_ar_p(0:42) #length = 42
predict_ar_p.head()
```

Out[125]: 9 0.146767
10 -0.022745
11 0.217014
12 -0.137244
13 0.136466
dtype: float64

iii. plot the predicted results and the corresponding seasonally differenced time-series

```
In [126]: #Plot the Predicted Results and the Corresponding Seasonally Differenced Time Series
plt.figure(figsize=(16, 3))
# Predict values
plt.grid()
plt.plot(predict_ar_p,label = 'Predicted Time Series', color = 'maroon')
# Seasonally Differenced Time Series
plt.plot(seas_best, label = 'Seasonally Differenced Time Series', color = 'gold')
plt.title('Predicted results and the corresponding Seasonally Differenced Time Series')
plt.xlabel('Time (Quarter)')
plt.legend()
plt.show()
```



iv. Calculate the Mean Absolute Error (MAE) by comparing the predicted results with the seasonally differenced data.

```
In [127]: #Finding MAE
maeq2 = mean_absolute_error(seas_best[5:], predict_ar_p)
maeq2
```

Out[127]: 0.13139640993766555

Mean Absolute Error is 0.13139640993766555

```
In [128]: mseq2 = mean_squared_error(seas_best[5:], predict_ar_p)
mseq2
```

Out[128]: 0.028569077607855303

Mean Squared Error is 0.028569077607855303

```
In [129]: sseqq2 = np.sum((seas_best[5:] - predict_ar_p)**2)
sseqq2
```

Out[129]: 1.1999012595299228

SSE is 1.1999012595299228

In []: