

In [ ]:

```

1  #####
2  #                               Session 32 -Assignment
3  #####
4  #Problem Statement: Students have to make ARIMA model over shampoo sales
5  #                               data and check the MSE between predicted and actual value
6
7
8  # importing the packages and load the data
9  import pandas as pd
10 import numpy as np
11
12 import matplotlib.pyplot as plt
13 % matplotlib inline
14 import statsmodels.api as sm
15 from statsmodels.tsa.arima_model import ARIMA
16 from statsmodels.tsa.stattools import acf, pacf
17 from sklearn.metrics import mean_squared_error
18
19 print("Shampoo Sales Data\n", '-'*80)
20 shmpData = pd.read_csv('shampoo-sales.csv', header=0, names = ['Date', 'Shampoo-sales'] )
21 shmpData['FDate'] = '190'+shmpData['Date']
22 shmpData.drop(['Date'], axis=1, inplace=True)
23 print(shmpData.head())
24

```

In [ ]:

```

1  print(" Analysing the Data \n", '-'*80)
2  print(shmpData.info())
3  print('-'*80)
4  print(shmpData.describe())
5  print('-'*80)
6  print(shmpData.isnull().sum())
7  print('-'*80)
8  shmpData.dropna(inplace=True)
9  print(shmpData.isnull().sum())
10 print('-'*80)
11
12 print(shmpData.shape)

```

```
In [ ]: 1 print(" Making the Date as the Index for the DataFrame\n", '- '*80)
        2 shmpData['FDate'] = pd.to_datetime(shmpData['FDate'])
        3 shmpData.set_index('FDate', inplace=True)
        4 print(shmpData.head())
        5
        6
```

```
In [ ]: 1 #Step1 : Visualising the Data
        2
        3 %matplotlib inline
        4 shmpData.plot(figsize=(6,4))
        5 plt.xlabel("Date")
        6 plt.ylabel("Shampoo Sales")
        7
        8 print(" We Observe from the plot , that there is upward trend ")
```

```
In [ ]: 1 # We observe from the graph above that there is a upward trend in the shampoo sale
2 # we need to check for stationarity
3
4 # Defining the the rolling mean and std deviation and Dickey_fuller to see the plot is stationary
5 from statsmodels.tsa.stattools import adfuller
6
7 # get_mean_var function gives the mean and the variance for the values
8 def get_mean_var(series, no_of_samples):
9     split_size = int(len(series) / no_of_samples)
10    start = 0
11    for i in range(no_of_samples):
12        sample_series = series[i*split_size:(i+1)*split_size]
13        #print(sample_series, "\n")
14        print('Mean= %.2f, Variance= %.2f' % (sample_series.mean(), sample_series.var()))
15
16
17 def plot_rolling_statistics(timeseries):
18     #Determining rolling statistics
19     rolmean = timeseries.rolling(window=12).mean()
20     rolstd = timeseries.rolling(window=12).std()
21
22     #Plot rolling statistics:
23     plt.plot(timeseries, color='blue', label='Original')
24     plt.plot(rolmean, color='red', label='Rolling Mean')
25     plt.plot(rolstd, color='black', label = 'Rolling Std')
26
27     plt.legend(loc='best')
28     plt.title('Rolling Mean & Standard Deviation')
29     plt.figure(figsize=(20,40))
30     plt.show()
31
32 # to perform Dickey fuller test
33
34 def dickey_fuller_test(timeseries):
35     #Perform Dickey-Fuller test:
36     print('Results of Dickey-Fuller Test:')
37     dftest = adfuller(timeseries['Shampoo-sales'], autolag='AIC')
38     #https://www.statsmodels.org/dev/generated/statsmodels.tsa.stattools.adfuller.html
39
40     dfoutput = pd.Series(dftest[0:4], index=['Test Statistic', 'p-value', '#Lags Used', 'Number of Observations Used'])
41
```

```
42     for key,value in dfctest[4].items():
43         dfoutput['Critical Value (%s)'%key] = value
44     print(dfoutput)
45
46
47
48
49
50 def test_stationarity(timeseries):
51     get_mean_var(timeseries,10)
52     plot_rolling_statistics(timeseries)
53     dickey_fuller_test(timeseries)
54
55
```

```
In [ ]: 1  # To see if the indexed Data is stationary by calling the above functions
2  test_stationarity(shmpData)
3
4  print('-'*80, ''\n From the Stationarity Test we Observe
5  1. P Value is big, should be 0
6  2. T stats should be Less than Critical Val , bu it is opposite for all confidence
7  hence we reject null , saying it is not stationary
8
9  ''')
```

```
In [ ]: 1  # Step 2 : Stationarizing the Time Series
2  import numpy as np
3  # Estimating the trend
4  shmpData_logScale = np.log(shmpData)
5
6  print(" Plotting Log Scale for Shampoo Data\n", '-'*80)
7  plt.plot(shmpData_logScale)
8  print(shmpData_logScale.head())
9
```

```
In [ ]: 1 test_stationarity(shmpData_logScale)
2 print('-'*80,'\nDurbin Watson Test :',sm.stats.durbin_watson(shmpData_logScale))
3
4 print('-'*80, ''\n From the Stationarity Test we Observe
5 1. P Value is big, should be 0
6 2. T stats should be Less than Critical Val , bu it is opposite for all confidence
7 hence we reject null , saying it is not stationary
8
9 ''')
```

```
In [ ]: 1 # DataSet with moving average
2 rolmean1 = shmpData_logScale.rolling(window=12).mean()
3 #print(rolmean)
4 shmpData_logScaleMinusMovingAvg = shmpData_logScale - rolmean1
5
6 # we remove the null values
7 shmpData_logScaleMinusMovingAvg.dropna(inplace=True)
8 shmpData_logScaleMinusMovingAvg.head()
```

```
In [ ]: 1 test_stationarity(shmpData_logScaleMinusMovingAvg)
2 print('-'*80,'\nDurbin Watson Test :',sm.stats.durbin_watson(shmpData_logScaleMinusMovingAvg))
3
4 print('-'*80, ''\n From the Stationarity Test we Observe
5 1. P Value is big, should be 0
6 2. T stats should be Less than Critical Val , bu it is opposite for all confidence
7 hence we reject null , saying it is not stationary
8
9 ''')
```

```
In [ ]: 1 # We Try Shifting the series , trying the Difference
2 shmpData_logScaleDiffShifting = shmpData_logScaleMinusMovingAvg - shmpData_logScaleMinusMovingAvg.shift()
3
4
5 shmpData_logScaleDiffShifting.dropna(inplace=True)
6 plt.plot(shmpData_logScaleDiffShifting)
```

```
In [ ]: 1 test_stationarity(shmpData_logScaleDiffShifting)
2
3 print('-'*80,'\nDurbin Watson Test :',sm.stats.durbin_watson(shmpData_logScaleDiffShifting))
4 print('-'*80, '''\n From the Stationarity Test we Observe
5 1. P Value is equal 0
6 2. T stats is Less than Critical Val for 90% and 95% confidence, we accept the null hypothesis
7 that time series is Stationary
8
9 ''')
```

```
In [ ]: 1 # test stats is less than 10% and 5% confidence
2
3 from statsmodels.tsa.seasonal import seasonal_decompose
4 decomposition = seasonal_decompose(shmpData_logScaleDiffShifting)
5 trend = decomposition.trend
6 seasonal = decomposition.seasonal
7 residual = decomposition.resid
8
9 plt.subplot(411)
10 plt.plot(shmpData_logScaleDiffShifting,label="Original")
11 plt.legend(loc="best")
12
13 plt.subplot(412)
14 plt.plot(trend,label="Trend")
15 plt.legend(loc="best")
16
17 plt.subplot(413)
18 plt.plot(seasonal,label="seasonal")
19 plt.legend(loc="best")
20
21 plt.subplot(414)
22 plt.plot(residual,label="Residual")
23 plt.legend(loc="best")
24
25 plt.tight_layout()
26
```

```
In [ ]: 1 decomposeLogData = residual
        2 decomposeLogData.dropna(inplace=True)
        3 test_stationarity(decomposeLogData)
        4 #print(decomposeLogData)
```

```
In [ ]: 1 # acf and pacf
        2 from statsmodels.tsa.stattools import acf, pacf
        3 lag_acf = acf(shmpData_logScaleDiffShifting, nlags=20)
        4 lag_pacf = pacf(shmpData_logScaleDiffShifting, nlags=20, method="ols")
        5
        6 #acf- gives q value where it touches the top confidence level
        7 plt.subplot(121)
        8 plt.plot(lag_acf)
        9 plt.axhline(y=0, linestyle='--', color='gray')
       10 plt.axhline(y=-1.96/np.sqrt(len(shmpData_logScaleDiffShifting)), linestyle='--', color='gray')
       11 plt.axhline(y= 1.96/np.sqrt(len(shmpData_logScaleDiffShifting)), linestyle='--', color='gray')
       12 plt.title("AutoCorrelation")
       13
       14 #pacf -gives p value where it touches the top confidence level
       15 plt.subplot(122)
       16 plt.plot(lag_pacf)
       17 plt.axhline(y=0, linestyle='--', color='gray')
       18 plt.axhline(y=-1.96/np.sqrt(len(shmpData_logScaleDiffShifting)), linestyle='--', color='gray')
       19 plt.axhline(y= 1.96/np.sqrt(len(shmpData_logScaleDiffShifting)), linestyle='--', color='gray')
       20 plt.title("Partial AutoCorrelation")
       21 plt.tight_layout()
       22
       23
```

```
In [ ]: 1 %matplotlib inline
2 import statsmodels.api as sm
3 fig = plt.figure(figsize=(12,8))
4
5 ax1 = fig.add_subplot(211)
6 fig = sm.graphics.tsa.plot_acf(shmpData_logScaleDiffShifting.values.squeeze(), lags=20, ax=ax1)
7
8 ax2 = fig.add_subplot(212)
9 fig = sm.graphics.tsa.plot_pacf(shmpData_logScaleDiffShifting.values.squeeze(), lags=20, ax=ax2)
```

```
In [ ]: 1 # We observe that p and q here is 1
2 # Ar model
3
4 # check is it p=1 and q=1 as per the graph
5 # also when we using arima model , the sites mention to put the Logged Dataframe , but we have achieved
6 # stationarity to DiffShiting
7 # For Moving Average model we get Least Least RSS
8
9 model = ARIMA(shmpData_logScaleDiffShifting,order=(0,1,1))
10 #model = ARIMA(shmpData_logScale,order=(1,1,1))
11 results_ARIMA = model.fit()
12 plt.plot(shmpData_logScaleDiffShifting)
13 plt.plot(results_ARIMA.fittedvalues,color="red")
14 plt.title('RSS: %.4f' % sum(((results_ARIMA.fittedvalues-shmpData_logScaleDiffShifting['Shampoo-sales']).dropna()))**2)
15 print("Plotting ARIMA Model")
16
```



```

In [ ]: # fitting the Arima Model and Predictions
1
2
predictions_ARIMA_diff = pd.Series(results_ARIMA.fittedvalues,copy=True)
4
# taking cumulativeSum
5
predictions_ARIMA_diff.cumsum = predictions_ARIMA_diff.cumsum()
6
# print(predictions_ARIMA_diff.cumsum.head())
7
# print('- '*80)
8
9
10
# predictions for the fitted values
11
# print(shmpData_logScale.head())
12
predictions_ARIMA_log = pd.Series(shmpData_logScale['Shampoo-sales'].values,index= shmpData_logScale.index)
13
14
rolmean1.fillna(0,inplace=True)
15
# print(rolmean1['Shampoo-sales'].values)
16
# print(predictions_ARIMA_diff.cumsum)
17
18
# predictions_ARIMA_mean = predictions_ARIMA_log.add(rolmean1,fill_value=0)
19
predictions_ARIMA_log = predictions_ARIMA_log.add(predictions_ARIMA_diff.cumsum,fill_value=0)
20
# predictions_ARIMA_Avg = predictions_ARIMA_log.add(rolmean1['Shampoo-sales'].values)
21
# print(predictions_ARIMA_log.head())
22
# print(predictions_ARIMA_Avg.head())
23
24
25
predictions_ARIMA = np.exp(predictions_ARIMA_log)
26
print(" Original Shampoo Data")
27
print(shmpData.head())
28
print(" Predicted Shampoo Data ")
29
print(predictions_ARIMA.head())
30
31
plt.plot(shmpData)
32
plt.plot(predictions_ARIMA)
33
plt.title('RMSE: %.4f'% np.sqrt(((predictions_ARIMA.values- shmpData.values)**2).sum()/len(shmpData)))
34
35

```

```

In [ ]: 1 results_ARIMA.plot_predict(1,60)

```

In [ ]:

```
1
2 def mean_forecast_err(y, yhat):
3     return y.sub(yhat).mean()
4
5 def mean_absolute_err(y, yhat):
6     return np.mean((np.abs(y.sub(yhat).mean()) / yhat)) # or percent error = * 100
7
8
9
10 print("Mean ForeCast Error - MFE = ",mean_forecast_err(shmpData['Shampoo-sales'], predictions_ARIMA))
11 print("Mean Absolute Error - MAE = ",mean_absolute_err(shmpData['Shampoo-sales'], predictions_ARIMA))
12 print("Mean Square Error - MSE = ",mean_squared_error(shmpData['Shampoo-sales'], predictions_ARIMA))
13
```