

## Import neccessery libraries

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import statsmodels.api as sm
from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.tsa.holtwinters import SimpleExpSmoothing
from statsmodels.tsa.holtwinters import Holt
from statsmodels.tsa.holtwinters import ExponentialSmoothing
import statsmodels.graphics.tsaplots as tsa_plots
import statsmodels.tsa.statespace as tm_models
from datetime import datetime, time
import warnings
import itertools
import matplotlib.pyplot as plt
warnings.filterwarnings("ignore")
plt.style.use('fivethirtyeight')
import pandas as pd
import statsmodels.api as sm
import matplotlib
from pylab import rcParams
from statsmodels.tsa.arima_model import ARIMA
from matplotlib import pyplot
from sklearn.metrics import mean_squared_error
```

## Problem

**Forecast the airlines data set. Prepare a document for each model explaining how many dummy variables you have created and RMSE value for each model. Finally which model you will use for Forecasting**

## Import data

```
In [2]: airline = pd.read_excel('Airlines+Data.xlsx')
```

Out[2]:

	Month	Passengers
0	1995-01-01	112
1	1995-02-01	118
2	1995-03-01	132
3	1995-04-01	129
4	1995-05-01	121
...	...	...
91	2002-08-01	405
92	2002-09-01	355
93	2002-10-01	306

	Month	Passengers
94	2002-11-01	271
95	2002-12-01	306

In [3]:

In [14]:

Out[14]:

	Month	Passengers
0	1995-01-01	112
1	1995-02-01	118
2	1995-03-01	132
3	1995-04-01	129
4	1995-05-01	121

## Data understanding

In [15]:

Out[15]: (96, 2)

In [6]:

Out[6]: Month 0  
Passengers 0  
dtype: int64

In [7]:

Out[7]: Month datetime64[ns]  
Passengers int64  
dtype: object

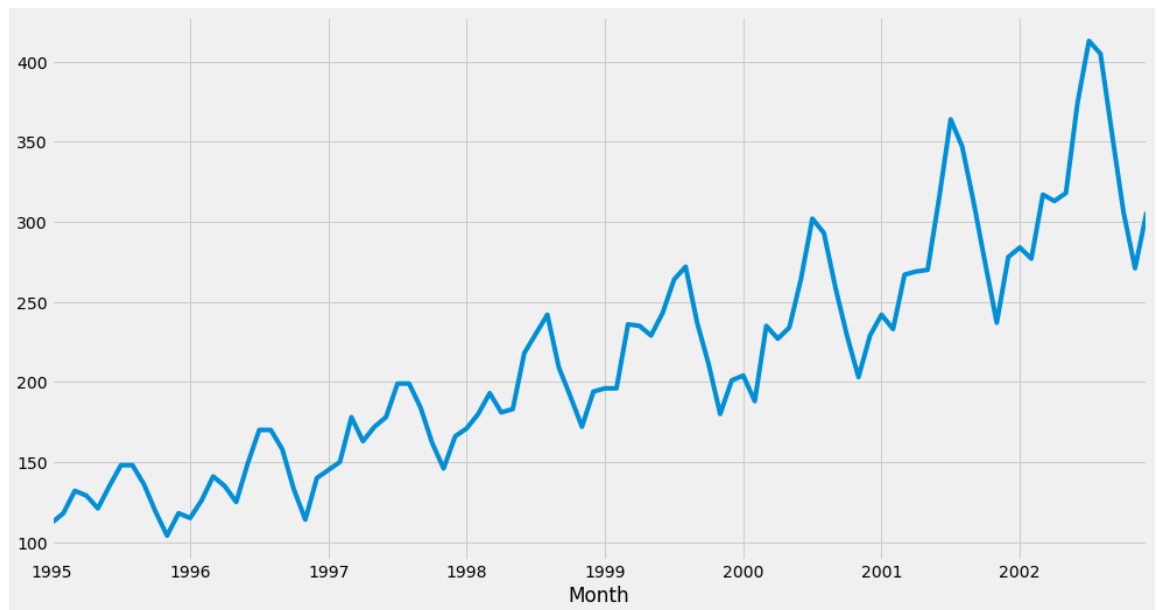
In [8]:

Out[8]:

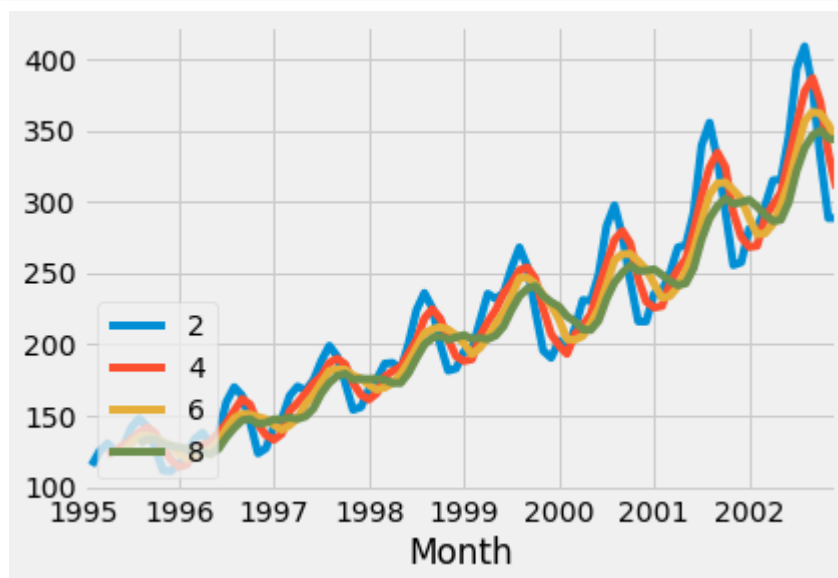
	count	mean	std	min	25%	50%	75%	max
Passengers	96.0	213.708333	71.918216	104.0	156.0	200.0	264.75	413.0

In [9]:

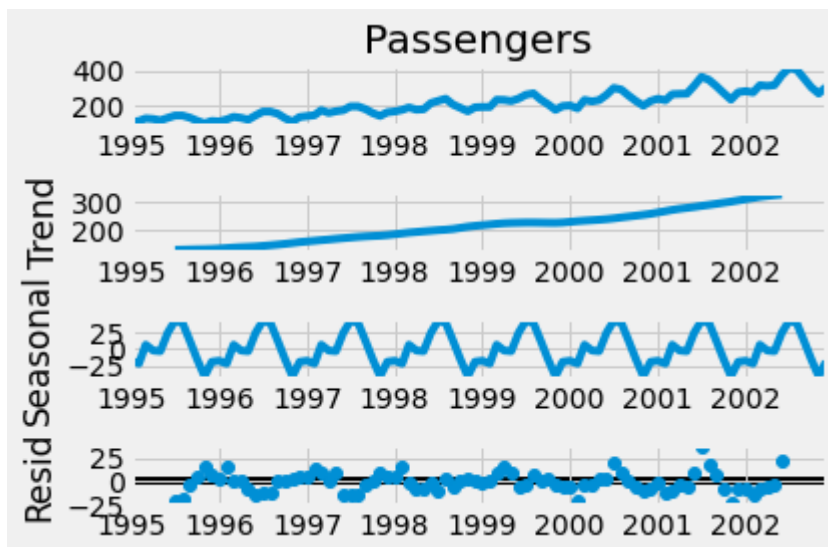
```
In [72]: airline_data_2['Passengers'].plot(figsize=(15,8))
```



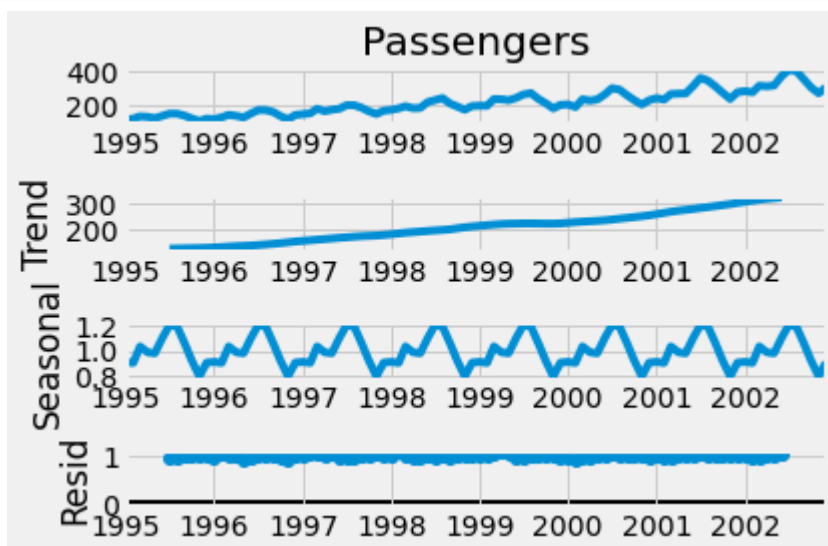
```
In [73]: for i in range(2,10,2):  
         airline_data_2['Passengers'].rolling(i).mean().plot(label=str(i))
```



```
In [74]: ts_add = seasonal_decompose(airline_data_2['Passengers'],model="additive")
fig = ts_add.plot()
```

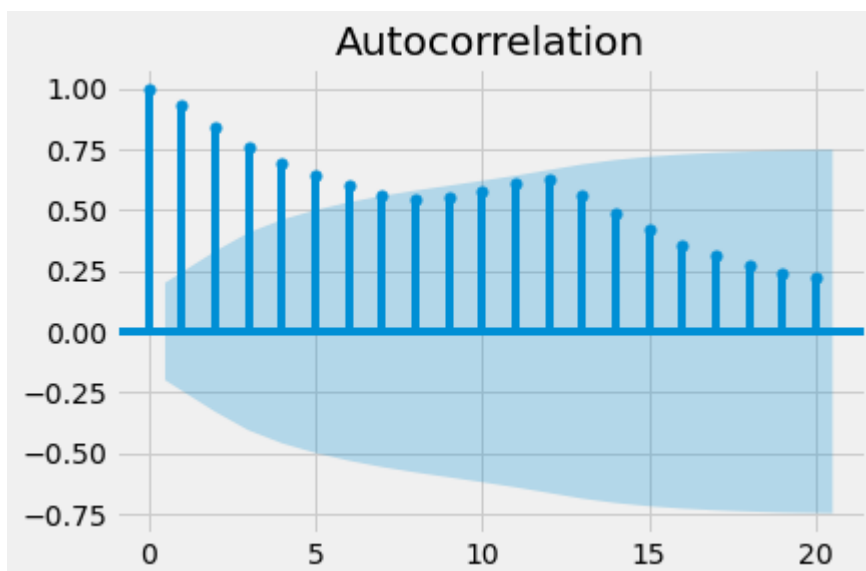


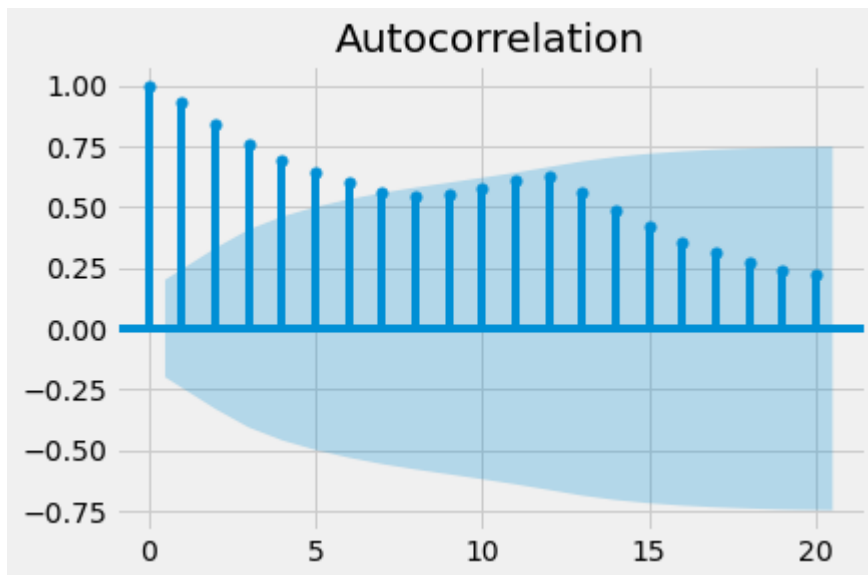
```
In [75]: ts_mul = seasonal_decompose(airline_data_2.Passengers,model="multiplicative")
fig = ts_mul.plot()
plt.show()
```



```
In [76]: plt.figure(figsize=(10,6))
plt.plot(airline_data_2['Passengers'].autocorr())
```

Out[76]:





## Building Time series forecasting with ARIMA

In [77]:

```
from statsmodels.tsa.arima.model import ARIMA
```

In [78]:

```
model = ARIMA(y, order=(1, 1, 1))
```

In [79]:

```
model = model.fit()
```

In [80]:

```
model.summary()
```

In [81]:

```
model.predict(20)
```

In [82]:

```
model.conf_int(20)
```

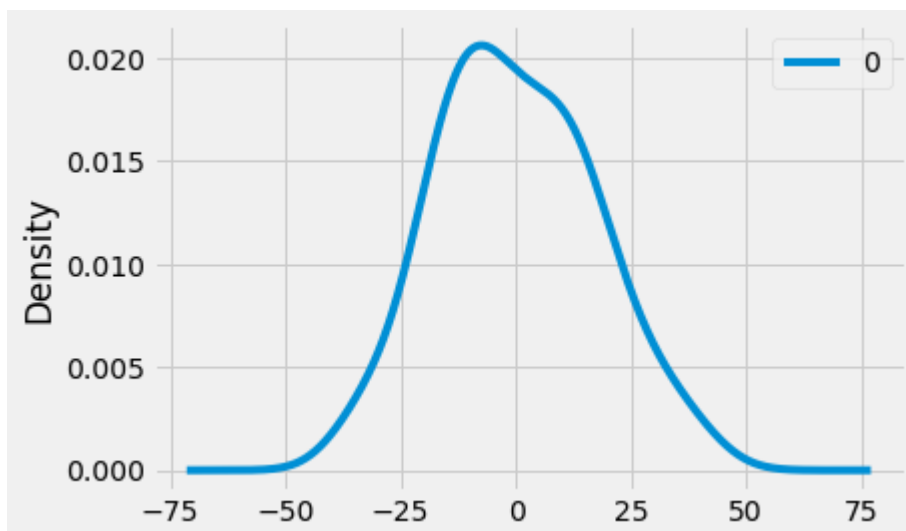
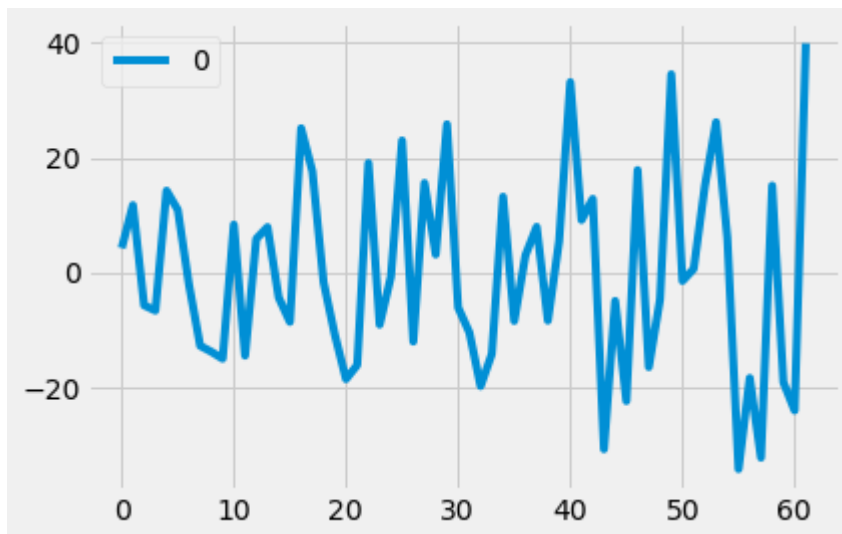
# ARIMA Model Results

```
=====
=====
Dep. Variable:          D.y    No. Observations:
62
Model:                ARIMA(5, 1, 0)    Log Likelihood
-262.909
Method:                css-mle    S.D. of innovations
16.748
Date:                  Sat, 26 Mar 2022    AIC
539.817
Time:                  11:01:36    BIC
554.707
Sample:                1    HQIC
545.663
```

```
=====
=====
               coef      std err          z      P>|z|      [0.025
0.975]
-----
-----
const          1.7497         1.477         1.185         0.236        -1.145
4.644
ar.L1.D.y       0.0905         0.134         0.677         0.498        -0.171
0.352
ar.L2.D.y      -0.2096         0.135        -1.549         0.121        -0.475
0.056
ar.L3.D.y      -0.0820         0.133        -0.622         0.532        -0.344
```

**This summarizes the coefficient values used as well as the skill of the fit on the on the in-sample observations**

```
In [83]: residuals = pd.DataFrame(model_fit.resid)
residuals.plot()
pyplot.show()
residuals.plot(kind='kde')
pyplot.show()
```



```
0
count    62.000000
mean      0.057356
std       16.895802
min      -34.303298
25%     -12.610648
50%      -1.589475
75%      12.565599
max       39.955366
```

**The plot of the residual errors suggests that there may still be some trend information not captured by the model**

The results show that there is no a bias in the prediction (a zero mean in the residuals)

## Rolling Forecast ARIMA Model

In [84]:

In [85]:

```
In [86]: for t in range(len(test)):
          model = ARIMA(history, order=(5,1,0))
          model_fit = model.fit(dispatch=0)
          output = model_fit.forecast()
          yhat = output[0]
          predictions.append(yhat)
          obs = test[t]
          history.append(obs)
```

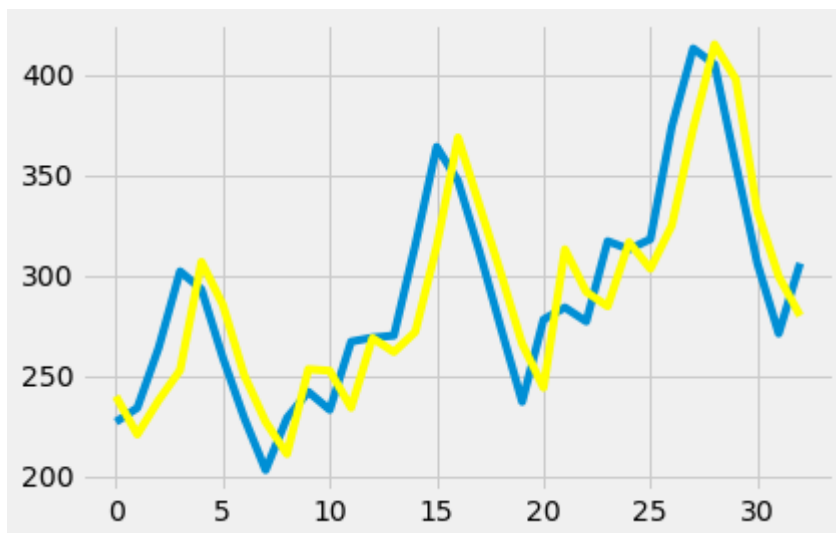
```
predicted=239.755179, expected=227.000000
predicted=220.737300, expected=234.000000
predicted=237.815008, expected=264.000000
predicted=252.750592, expected=302.000000
predicted=306.715794, expected=293.000000
predicted=285.374643, expected=259.000000
predicted=250.264004, expected=229.000000
predicted=227.093124, expected=203.000000
predicted=211.011455, expected=229.000000
predicted=253.260281, expected=242.000000
predicted=252.490682, expected=233.000000
predicted=234.042132, expected=267.000000
predicted=268.773632, expected=269.000000
predicted=261.782257, expected=270.000000
predicted=271.798054, expected=315.000000
predicted=314.422115, expected=364.000000
predicted=368.637742, expected=347.000000
predicted=334.957878, expected=312.000000
predicted=301.161818, expected=274.000000
predicted=265.936481, expected=237.000000
predicted=244.037181, expected=278.000000
predicted=312.961792, expected=284.000000
predicted=291.748167, expected=277.000000
predicted=284.551868, expected=317.000000
predicted=316.501195, expected=313.000000
predicted=303.218148, expected=318.000000
predicted=324.834619, expected=374.000000
predicted=373.140656, expected=413.000000
predicted=415.007200, expected=405.000000
predicted=397.508453, expected=355.000000
predicted=332.087112, expected=306.000000
predicted=299.452956, expected=271.000000
predicted=279.908341, expected=306.000000
```

```
In [87]: error = mean_squared_error(test, predictions)
```

Test MSE: 782.495



```
In [88]: pyplot.plot(test)
pyplot.plot(predictions, color='yellow')
```



**A line plot is created showing the expected values (blue) compared to the rolling forecast predictions (yellow). We can see the values show some trend and are in the correct scale**

## Comparing Multiple Models

```
In [51]: airline_data_2 = airline_data_1[['M', 'Passengers']]
```

```
In [52]: airline_data_2['M'] = pd.to_datetime(airline_data_2['M'], format='%m-%Y')
```

```
In [50]: airline_data_2.head()
```

```
Out[50]:
```

	Passengers	Month_1995-01-01 00:00:00	Month_1995-02-01 00:00:00	Month_1995-03-01 00:00:00	Month_1995-04-01 00:00:00	M
0	112	1	0	0	0	
1	118	0	1	0	0	
2	132	0	0	1	0	
3	129	0	0	0	1	
4	121	0	0	0	0	

5 rows × 97 columns

```
In [20]: airline_data_2 = airline_data_2[['M', 'Passengers']]
```

```
Out[20]: (96, 97)
```

```
In [21]: airline_data_2['M'] = pd.to_datetime(airline_data_2['M'], format='%m-%Y')
```

```
In [22]: airline_data_2['M'] = pd.to_datetime(airline_data_2['M'], format='%m-%Y')
```

```
In [23]: airline_data_3['t_sq'] = airline_data_3['t'] * airline_data_3['t']
```

```
In [33]: passengers1= np.log(airline_data_3['Passengers'])
```

```
In [34]:
```

```
In [35]:
```

```
Out[35]:
```

	Passengers	Month_1995-01-01 00:00:00	Month_1995-02-01 00:00:00	Month_1995-03-01 00:00:00	Month_1995-04-01 00:00:00	M
0	4.718499	1	0	0	0	
1	4.770685	0	1	0	0	
2	4.882802	0	0	1	0	
3	4.859812	0	0	0	1	
4	4.795791	0	0	0	0	

5 rows × 100 columns

```
In [36]:
```

```
In [37]: linear= smf.ols('Passengers ~ t',data=train1).fit()  
predlin=pd.Series(linear.predict(pd.DataFrame(test1['t'])))  
rmselin=np.sqrt((np.mean(np.array(test1['Passengers'])-np.array(predlin)  
rmselin
```

```
Out[37]: 0.011446730996560794
```

```
In [38]: quad=smf.ols('Passengers~t+t_sq',data=train1).fit()  
predquad=pd.Series(quad.predict(pd.DataFrame(test1[['t','t_sq']])))  
rmsequad=np.sqrt(np.mean((np.array(test1['Passengers'])-np.array(predquad)  
rmsequad
```

```
Out[38]: 0.17833090054825948
```

```
In [39]: expo=smf.ols('Passengers~t',data=train1).fit()  
predexp=pd.Series(expo.predict(pd.DataFrame(test1['t'])))  
rmseexpo=np.sqrt(np.mean((np.array(test1['Passengers'])-np.array(predexp)  
rmseexpo
```

```
Out[39]: 289.1236843586758
```

## Conclusion

```
In [66]: output = {'Model':pd.Series(['rmseexpo','rmselin','rmsequad']),  
                  'RMSE':pd.Series([rmseexpo,rmselin,rmsequad])}
```

```
In [67]: rmse=pd.DataFrame(output)
```

```
In [68]:
```

	Model	Values
0	rmseexpo	rmseexpo
1	rmselin	rmselin
2	rmsequad	rmsequad

In [ ]: