n [1]:	import Necessary Libraries import pandas as pd from matplotlib import pyplot as plt import seaborn as sns import statsmodels.api as sm from statsmodels.tsa.seasonal import seasonal_decompose import statsmodels.graphics.tsaplots as tsa_plots from statsmodels.tsa.arima_model import ARIMA from matplotlib import pyplot from sklearn.metrics import mean_squared_error import numpy as np import statsmodels.formula.api as smf
[77]: [81]:	Data Collection df=pd.read_excel('CocaCola_Sales_Rawdata.xlsx') cola=df.copy() cola.head()
[83]:	cola.shape (42, 2) cola.isna().sum() Quarter 0
[84]:	Sales 0 dtype: int64 cola.dtypes Quarter object Sales float64 dtype: object cola.describe().T
[86]:	count mean std min 25% 50% 75% max Sales 42.0 2994.353308 977.930896 1547.818996 2159.714247 2782.376999 3609.25 5253.0 import warnings warnings.filterwarnings('ignore') 100.0000000000000000000000000000000000
[87]: [88]: [89]: [89]:	<pre>temp = cola.Quarter.str.replace(r'(Q\d)_(\d+)', r'19\2-\1') cola['quater'] = pd.to_datetime(temp).dt.strftime('%b-%Y') cola.head() Quarter</pre>
	0 Q1_86 1734.827000 Jan-1986 1 Q2_86 2244.960999 Apr-1986 2 Q3_86 2533.804993 Jul-1986 3 Q4_86 2154.962997 Oct-1986 4 Q1_87 1547.818996 Jan-1987
[91]: [[92]: [[92]: _	<pre>cola = cola.drop(['Quarter'], axis=1) cola.reset_index(inplace=True) cola.head() index</pre>
[93]:	0 1734.827000 Jan-1986 1 1 2244.960999 Apr-1986 2 2 2533.804993 Jul-1986 3 3 2154.962997 Oct-1986 4 4 1547.818996 Jan-1987 cola['Sales'].plot(figsize=(15, 6))
	plt.show() 5000 - 4500 -
	3500 - 3000 - 2500 - 2000 -
[94]:	for i in range(2,10,2): cola["Sales"].rolling(i).mean().plot(label=str(i)) plt.legend(loc=3) <matplotlib.legend.legend 0x1fda6651f40="" at=""></matplotlib.legend.legend>
	4500 - 4000 - 3500 - 2500 -
[95]:	ts_add = seasonal_decompose(cola.Sales,model="additive", period=2) fig = ts_add.plot() plt.show() Sales
	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$
[96]:	ts_mul = seasonal_decompose(cola.Sales, model="multiplicative", period=2) fig = ts_mul.plot() plt.show()
	Sales 5000 2500 5000
[97]: [97]:	tsa_plots.plot_acf(cola.Sales) Autocorrelation
	1.00 -
	-0.75 - 0.0 2.5 5.0 7.5 10.0 12.5 15.0 17.5 Autocorrelation
	0.25 - 0.00 -0.25 - 0.50 -0.75 - 0.00 12.5 15.0 17.5
[98]: [99]: [100	Building Time Series Forecasting X = cola['Sales'].values size = int(len(X) * 0.66) train, test = X[0:size], X[size:len(X)]
[101	<pre>model = ARIMA(train, order=(5,1,0)) model_fit = model.fit(disp=0) print(model_fit.summary())</pre>
	ARIMA Model Results ===================================
	const 41.8434 26.509 1.578 0.114 -10.113 93.799 ar.L1.D.y -0.1479 0.195 -0.758 0.448 -0.530 0.234 ar.L2.D.y -0.3127 0.157 -1.996 0.046 -0.620 -0.006 ar.L3.D.y -0.1881 0.173 -1.090 0.276 -0.526 0.150 ar.L4.D.y 0.6222 0.167 3.716 0.000 0.294 0.950 ar.L5.D.y -0.1766 0.220 -0.804 0.422 -0.607 0.254 Roots
	AR.3
	print(residuals.describe()) 400 -
	-200
	0.0015 - 0.0005 - 0.0005 - 0.0000 - 0.0
	count 26.000000 mean 31.325878 std 202.029702 min -438.904389 25% -58.603106 50% -9.190229 75% 200.236018 max 468.290585 Forecast ARIMA Model
[105	<pre>history = [x for x in train] predictions = list() for t in range(len(test)): model = ARIMA(history, order=(5,1,0)) model_fit = model.fit(disp=0)</pre>
	<pre>output = model_fit.forecast() yhat = output[0] predictions.append(yhat) obs = test[t] history.append(obs) print('predicted=%f, expected=%f' % (yhat, obs)) predicted=3135.586029, expected=3243.859993 predicted=3188.847068, expected=3056.000000 predicted=3734.224502, expected=3899.000000</pre>
	predicted=3782.620891, expected=3629.000000 predicted=3355.125969, expected=3373.000000 predicted=3297.218120, expected=3352.000000 predicted=4112.813891, expected=4342.000000 predicted=43961.043678, expected=4461.000000 predicted=4130.787225, expected=4017.000000 predicted=4970.516924, expected=4936.000000 predicted=4970.516924, expected=4936.000000 predicted=4384.040534, expected=4333.000000 predicted=4207.687405, expected=4194.000000 predicted=5261.673040, expected=4194.000000 expected=5253.000000
[108	error = mean_squared_error(test, predictions) print('Test MSE: %.3f' % error) Test MSE: 31525.273 pyplot.plot(test) pyplot.plot(predictions, color='red') pyplot.show()
	5000 - 4500 - 4000 -
[116	3500 3000 Comparing Multiple Models cola2 = pd.get_dummies(df, columns = ['Quarter'])
[117	cola2.columns = ['Sales','Q1','Q1','Q1','Q1','Q1','Q1','Q1','Q
	1 2244.960999 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
[120	<pre>t= np.arange(1,43) cola2['t'] = t cola2['t_sq'] = cola2['t']*cola2['t'] log_Sales=np.log(cola2['Sales'])</pre>
	cola2['log_Sales']=log_Sales']=log_Sales']=log_Sales']=log_Sales']=log_Sales
	1 2244.960999 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
[128	<pre>train1, test1 = np.split(cola2, [int(.67 *len(cola2))]) linear= smf.ols('Sales ~ t', data=train1).fit() predlin=pd.Series(linear.predict(pd.DataFrame(test1['t']))) rmselin=np.sqrt((np.mean(np.array(test1['Sales'])-np.array(predlin))**2)) rmselin 580.1224130918641</pre>
[129	<pre>quad=smf.ols('Sales~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['Sales'])-np.array(predquad))**2)) rmsequad 783.7297975037103 expo=smf.ols('log_Sales~t', data=train1).fit() predexp=pd.Series(expo.predict(pd.DataFrame(test1['t'])))</pre>
[130	
[132	<pre>rmseadd 1869.7188209186947 addlinear= smf.ols('Sales~t+Q1+Q2+Q3+Q4', data=train1).fit() predaddlinear=pd.Series(addlinear.predict(pd.DataFrame(test1[['t','Q1','Q2','Q3','Q4']]))) rmseaddlinear=np.sqrt(np.mean((np.array(test1['Sales'])-np.array(predaddlinear))**2)) rmseaddlinear</pre> 596.1526282372472
[133	<pre>addquad=smf.ols('Sales~t+t_sq+Q1+Q2+Q3+Q4',data=train1).fit() predaddquad=pd.Series(addquad.predict(pd.DataFrame(test1[['t','t_sq','Q1','Q2','Q3','Q4']]))) rmseaddquad=np.sqrt(np.mean((np.array(test1['Sales'])-np.array(predaddquad))**2)) rmseaddquad</pre> 412.1144436053775
4	<pre>mulsea=smf.ols('log_Sales-Q1+Q2+Q3+Q4', data=train1).fit() predmul= pd.Series(mulsea.predict(pd.DataFrame(test1[['Q1', 'Q2', 'Q3', 'Q4']]))) rmsemul= np.sqrt(np.mean((np.array(test1['Sales'])-np.array(np.exp(predmul)))**2)) rmsemul 2374.9194407954374 mullin= smf.ols('log_Sales~t+Q1+Q2+Q3+Q4', data=train1).fit() predmullin= pd.Series(mullin.predict(pd.DataFrame(test1[['t', 'Q1', 'Q2', 'Q3', 'Q4']])))</pre>
[134	rmsemulin=np.sqrt(np.mean((np.array(test1['Sales'])-np.array(np.exp(predmullin)))**2))
[135	<pre>rmsemulin=np.sqrt(np.mean((np.array(test1['Sales'])-np.array(np.exp(predmullin)))**2)) rmsemulin 5359.687911932085 mul_quad= smf.ols('log_Sales-t+t_sq+Q1+Q2+Q3+Q4', data=train1).fit() pred_mul_quad= pd.Series(mul_quad.predict(test1[['t','t_sq','Q1','Q2','Q3','Q4']])) rmse_mul_quad=np.sqrt(np.mean((np.array(test1['Sales'])-np.array(np.exp(pred_mul_quad)))**2)) rmse_mul_quad</pre>
[135 [135 [136 [137	<pre>rmsemulin=np.sqrt(np.mean((np.array(test1['Sales'])-np.array(np.exp(predmullin)))**2)) rmsemulin 5359.687911932085 mul_quad= smf.ols('log_Sales-t+t_sq+Q1+Q2+Q3+Q4', data=train1).fit() pred_mul_quad= pd.Series(mul_quad.predict(test1[['t','t_sq','Q1','Q2','Q3','Q4']])) rmse_mul_quad=np.sqrt(np.mean((np.array(test1['Sales'])-np.array(np.exp(pred_mul_quad)))**2)) rmse_mul_quad output = {'Model':pd.Series(['rmse_mul_quad', 'rmseadd', 'rmseaddlinear', 'rmseaddquad', 'rmseexpo', 'rmselin', 'rmsemul', 'rmsemulin', 'rmsequad']),</pre>
[135 [135 [136 [137 [138 [139	<pre>rmsemulin=np.sqrt(np.mean((np.array(test1['Sales'])-np.array(np.exp(predmullin)))**2)) rmsemul quad= smf.ols('log_Sales-t+t_sq+0l+02+03+04',data=train1).fit() pred_mul_quad= pd.Series(mul_quad.predict(test1[['t', 't_sa','01','02','03','04']])) rmse_mul_quad=np.sqrt(np.mean((np.array(test1['Sales'])-np.array(np.exp(pred_mul_quad)))**2)) rmse_mul_quad=np.sqrt(np.mean((np.array(test1['Sales'])-np.array(np.exp(pred_mul_quad)))**2)) rmse_mul_quad 3630.5619467347524 output = {'Model':pd.Series(['rmse_mul_quad', 'rmseadd', 'rmseaddlinear', 'rmseaddquad', 'rmseexpo', 'rmselin', 'rmsemul', 'rmsemulin', 'rmsequad']),</pre>
[135 [135 [136 [137 [138 [139	<pre>rmsemulin=np.sqrt(np.mean((np.array(testi['Sales'])-np.array(np.exp(predmullin)))**2)) rmsemulin 5359.687911932085 mul_quad= smf.ols('log_Sales-t+t_sq+01+02+03+04', data=train1).fit() pred_mul_quad= pd.Series(mul_quad.predict(testi['t','t_sq','Q1','Q2','Q3','Q4']))) rmse_mul_quad= pd.Series(mul_quad.predict(testi['t','t_sq','Q1','Q2','Q3','Q4']))) rmse_mul_quad 3630.5619467347524 output = {'Model':pd.Series(['rmse_mul_quad', 'rmseadd', 'rmseaddlinear', 'rmseaddquad', 'rmseexpo', 'rmselin', 'rmsemul', 'rmsemulin', 'rmsequad']),</pre>