

# Time Series Forecasting

## PROJECT - 7

### REPORT – Part 1

Piyush Kumar Singh  
PGP – DSBA Online  
May-21 Batch

Date: 16/01/2022

## Table of Contents

• List of Figures .....	3
• List of Tables.....	4
• Dataset 1 – Sparkling.csv .....	5
1.1 ) Read the data as an appropriate Time Series data and plot the data.....	5
1.2 ) Perform appropriate Exploratory Data Analysis to understand the data and also perform decomposition .....	6
1.3 ) Split the data into training and test. The test data should start in 1991 .....	11
1.4 ) Build various exponential smoothing models on the training data and evaluate the model using RMSE on the test data. Other models such as regression, naïve forecast models, simple average models etc. should also be built on the training data and check the performance on the test data using RMSE.....	13
1.5 ) Check for the stationarity of the data on which the model is being built on using appropriate statistical tests and also mention the hypothesis for the statistical test. If the data is found to be non-stationary, take appropriate steps to make it stationary. Check the new data for stationarity and comment. ....	15
1.6 ) Build an automated version of the ARIMA/SARIMA model in which the parameters are selected using the lowest Akaike Information Criteria (AIC) on the training data and evaluate this model on the test data using RMSE .....	17
1.7 ) Build ARIMA/SARIMA models based on the cut-off points of ACF and PACF on the training data and evaluate this model on the test data using RMSE. ....	20
1.8 ) Build a table with all the models built along with their corresponding parameters and the respective RMSE values on the test data.....	24
1.9 ) Based on the model-building exercise, build the most optimum model(s) on the complete data and predict 12 months into the future with appropriate confidence intervals/bands .....	25
1.10 ) Comment on the model thus built and report your findings and suggest the measures that the company should be taking for future sales.....	26

## List of Figures

Fig.1 – Line plot of rose Time Series along with mean and median of our Time Series .....	5
Fig.2 – Yearly boxplot of sparkling sales .....	6
Fig.3 – Monthly boxplot of sparkling sales .....	7
Fig.4 – Monthly plot showing mean and variation of units sold .....	7
Fig.5 – Line chart of sales across years .....	8
Fig.6 – Empirical cumulative distribution graph .....	9
Fig.7 – Average sales and sales percentage change across years .....	9
Fig.8 – Additive model decomposition .....	10
Fig.9 – Multiplicative model decomposition .....	10
Fig.10 – Original Time Series Vs Time Series with decomposed component .....	11
Fig.11 – Train and Test set line plot .....	12
Fig.12 – Prediction graph of various models w.r.t to test set .....	15
Fig.13 – Time Series of difference of order 1 .....	16
Fig.14 – Diagnostic plot of ARIMA(2,1,2) .....	18
Fig.15 – Diagnostic plot of SARIMA(3,1,2)(3,0,1,12) .....	19
Fig.16 – ACF and PACF plot of differenced time series .....	20
Fig.17 – Diagnostic plot of ARIMA(0,1,0) .....	21
Fig.18 – Diagnostic plot of SARIMA(0,1,0)(0,0,0,12) .....	23
Fig.19 – Line plot of Time Series Vs fitted values by model vs future prediction for 12 months .....	25
Fig.20 – Line plot of actual data Vs predicated data with 95% confidence interval .....	26
Fig.21 – Line plot of 12 months future prediction with 95% confidence interval .....	26

## List of Tables

Table 1. Summary of Time Series Sparkling Dataset .....	6
Table 2. Monthly sales across the years.....	8
Table 3. ARIMA Automated Summary results.....	17
Table 4. SARIMA Automated Summary results .....	19
Table 5. ARIMA Manual Summary results .....	21
Table 6. SARIMA Manual Summary results .....	22
Table 7. Model comparison table using test RMSE.....	24
Table 8. Sales Prediction values for 12 months in future.....	25

## Dataset: Sparkling.csv

For this particular assignment, the data of different types of wine sales in the 20th century is to be analysed. Both of these data are from the same company but of different wines. As an analyst in the ABC Estate Wines, you are tasked to analyse and forecast Wine Sales in the 20th century.

### 1.1 Read the data as an appropriate Time Series data and plot the data.

- Monthly sales of 'sparkling' wine from period Jan – 1980 to July – 1995 is provided in the sparkling.csv file.
- The given data files is read and date range is inserted and the YearMonth column is to date-range and set as index to create a time series data having one column of 'Sparkling' showing sales value.

Sparkling	
YearMonth	
1980-01-31	1686
1980-02-29	1591
1980-03-31	2304
1980-04-30	1712
1980-05-31	1471

Head

Sparkling	
YearMonth	
1995-03-31	1897
1995-04-30	1862
1995-05-31	1670
1995-06-30	1688
1995-07-31	2031

Tail

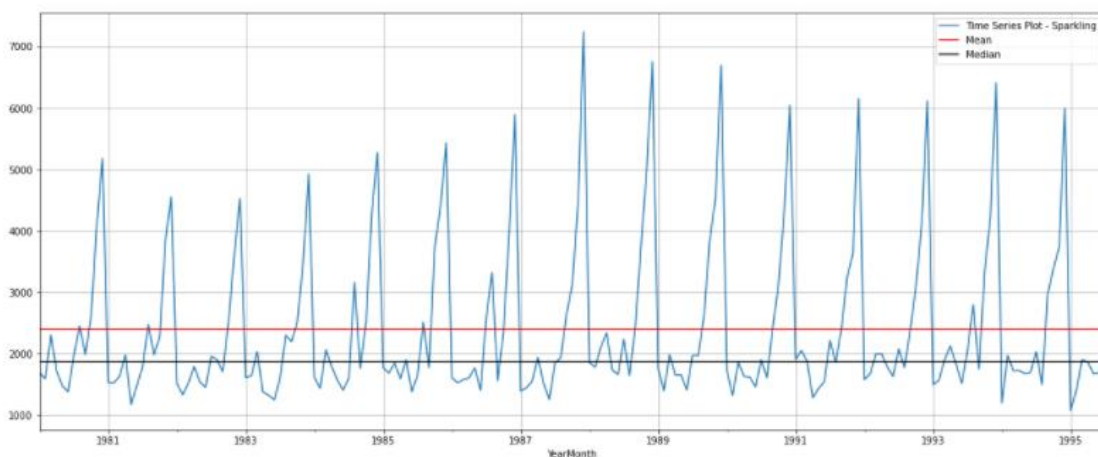


Fig.1

- Sparkling doesn't show any consistent trend and has both upward and downward slope during the mentioned time period.
- Sparkling wine sale has been consistent during the time period analysed. Seasonality is present on yearly basis.

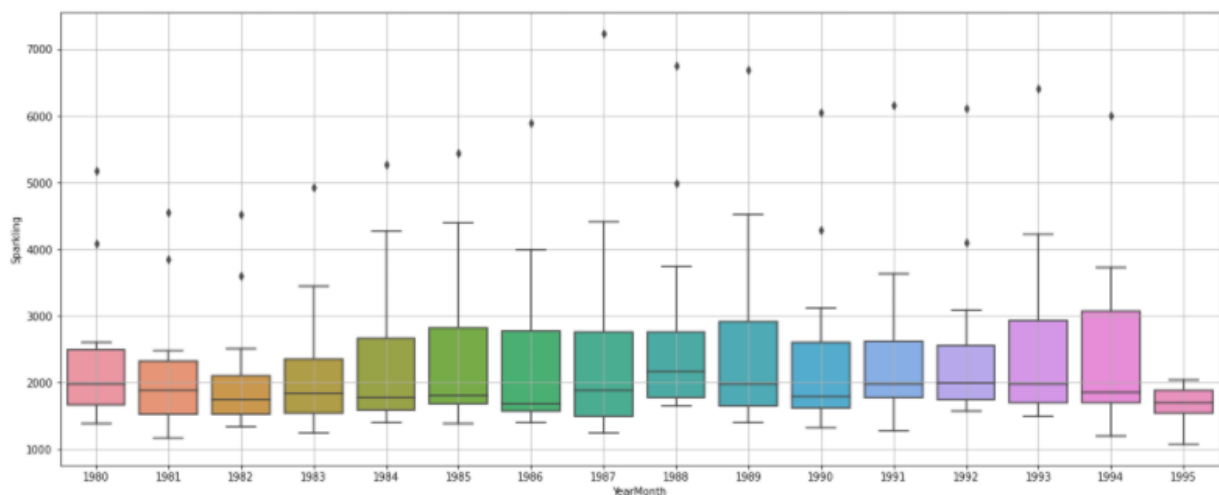
**1.2** Perform appropriate Exploratory Data Analysis to understand the data and also perform decomposition Read the data as an appropriate Time Series data and plot the data.

- Numbers of records in the dataset is 187.
- There are no null or missing values in our time series.

Sparkling	
count	187.000
mean	2402.417
std	1295.112
min	1070.000
25%	1605.000
50%	1874.000
75%	2549.000
max	7242.000

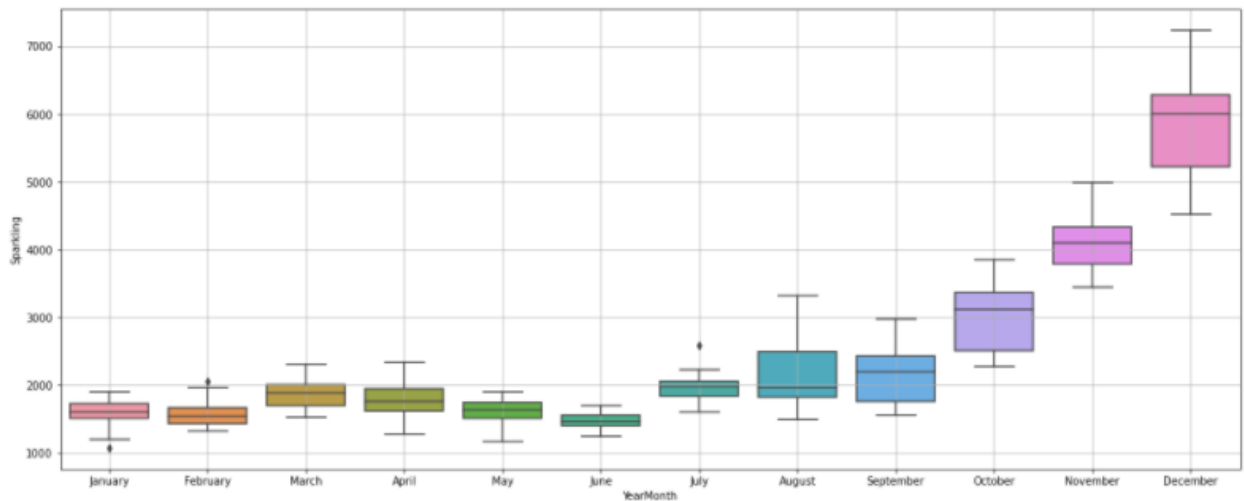
**Table.1**

- The descriptive summary of the data shows that on an average 2402 units of sparkling were sold each month on the given time period. 50 % of the sales have more than 1874 sales every month.
- Maximum sale in a month was 7242 and minimum sale in a month was 1070.



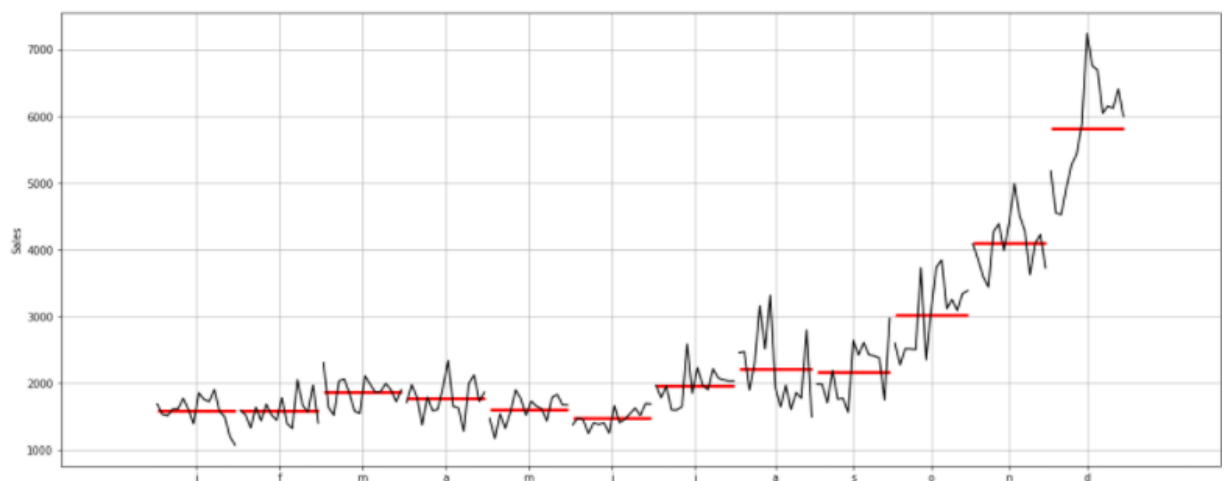
**Fig.2**

- The yearly boxplot shows that the sales has been more or less been the same throughout the years with sales close to approx. 2000 units across all the years.
- The outliers in the yearly-boxplot most probably represents seasonal sales during the seasonal months.



**Fig.3**

- The monthly boxplot show sales within different months spread across various years.
- The monthly boxplot shows clear seasonality during the months of September, October, November and December.
- The highest such numbers are being recorded in the month of December across various years.



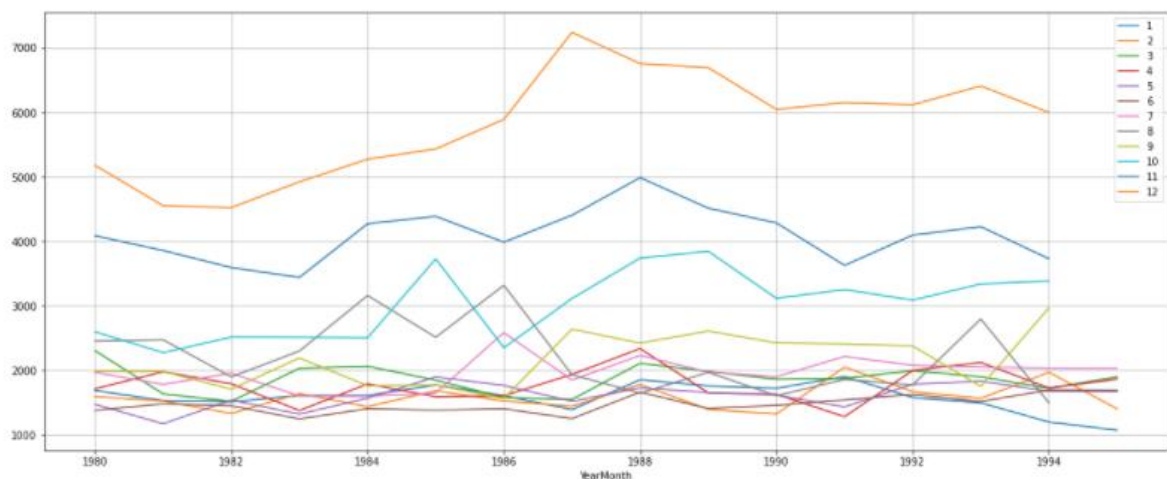
**Fig.4**

- This monthly plot shows mean and variation of units sold each month across various years.
- Sales in the seasonal months show higher variation.
- The red line is the mean sales value of each month across various years.
- Sales from Jan - June have been more or less the same across all the years and increases from July-Dec.

YearMonth	1	2	3	4	5	6	7	8	9	10	11	12
YearMonth												
1980	1686.0	1591.0	2304.0	1712.0	1471.0	1377.0	1966.0	2453.0	1984.0	2596.0	4087.0	5179.0
1981	1530.0	1523.0	1633.0	1976.0	1170.0	1480.0	1781.0	2472.0	1981.0	2273.0	3857.0	4551.0
1982	1510.0	1329.0	1518.0	1790.0	1537.0	1449.0	1954.0	1897.0	1706.0	2514.0	3593.0	4524.0
1983	1609.0	1638.0	2030.0	1375.0	1320.0	1245.0	1600.0	2298.0	2191.0	2511.0	3440.0	4923.0
1984	1609.0	1435.0	2061.0	1789.0	1567.0	1404.0	1597.0	3159.0	1759.0	2504.0	4273.0	5274.0
1985	1771.0	1682.0	1846.0	1589.0	1896.0	1379.0	1645.0	2512.0	1771.0	3727.0	4388.0	5434.0
1986	1606.0	1523.0	1577.0	1605.0	1765.0	1403.0	2584.0	3318.0	1562.0	2349.0	3987.0	5891.0
1987	1389.0	1442.0	1548.0	1935.0	1518.0	1250.0	1847.0	1930.0	2638.0	3114.0	4405.0	7242.0
1988	1853.0	1779.0	2108.0	2336.0	1728.0	1661.0	2230.0	1645.0	2421.0	3740.0	4988.0	6757.0
1989	1757.0	1394.0	1982.0	1650.0	1654.0	1406.0	1971.0	1968.0	2608.0	3845.0	4514.0	6694.0
1990	1720.0	1321.0	1859.0	1628.0	1615.0	1457.0	1899.0	1605.0	2424.0	3116.0	4286.0	6047.0
1991	1902.0	2049.0	1874.0	1279.0	1432.0	1540.0	2214.0	1857.0	2408.0	3252.0	3627.0	6153.0
1992	1577.0	1667.0	1993.0	1997.0	1783.0	1625.0	2076.0	1773.0	2377.0	3088.0	4096.0	6119.0
1993	1494.0	1564.0	1898.0	2121.0	1831.0	1515.0	2048.0	2795.0	1749.0	3339.0	4227.0	6410.0
1994	1197.0	1968.0	1720.0	1725.0	1674.0	1693.0	2031.0	1495.0	2968.0	3385.0	3729.0	5999.0
1995	1070.0	1402.0	1897.0	1862.0	1670.0	1688.0	2031.0	NaN	NaN	NaN	NaN	NaN

**Table.2**

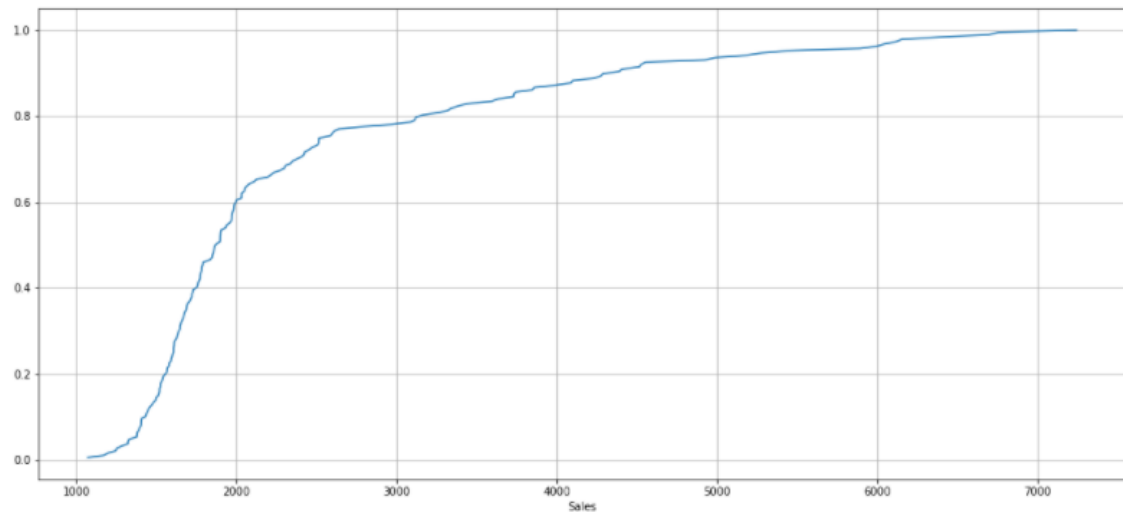
- Table of monthly sales across all the years.
- Highest sale in December month was in 1987 with units sold 7242.
- Least sale in December month was in 1982 with units sold 4524.



**Fig.5**

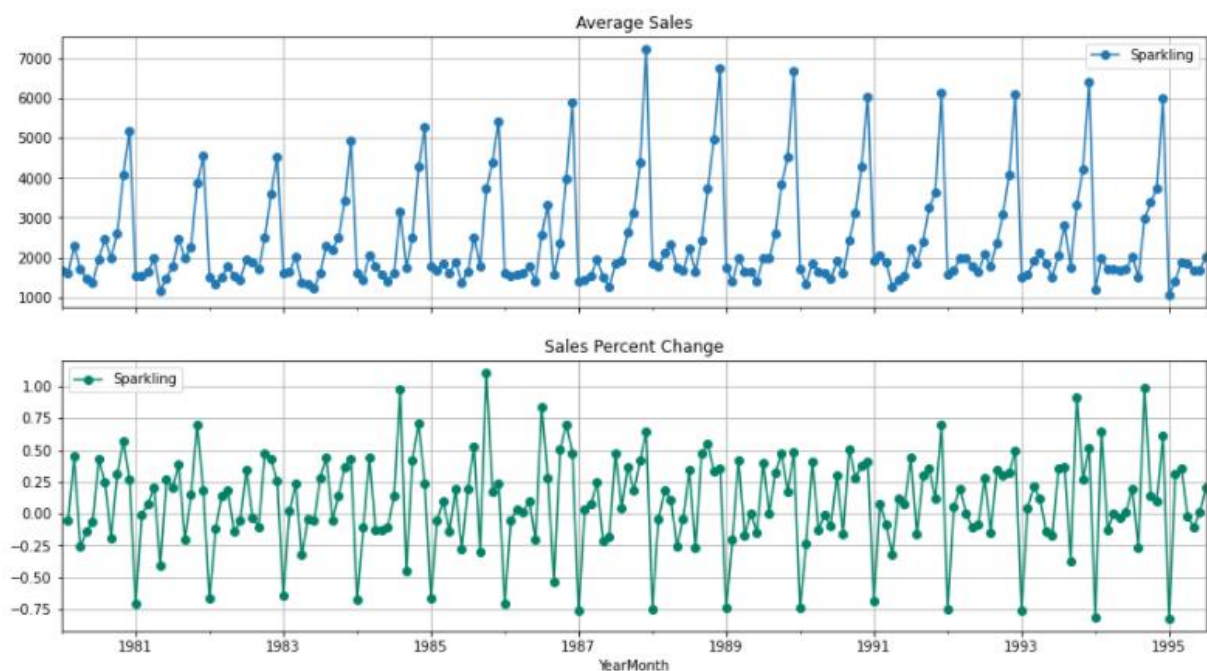
- Line chart of sales across various years. We can see December month has the highest number of sales of sparkling wine.
- Highest units sold in December was in 1987. Unit's sales crossed 7000 units in this year.





**Fig.6**

- This is an empirical cumulative distribution graph. This graph tells us what percentage of data points refer to what number of sales.
- 80% of the month have at least 3000 units' sales of sparkling wine.



**Fig.7**

- The above two graphs tell us the Average 'Sales' and the Percentage change of 'Sales' with respect to the time.
- Average sales shows a slight increase in units sold in later years crossing 6000 units' threshold sales during the seasonal period.

## Time Series Decomposition:

- Additive Model

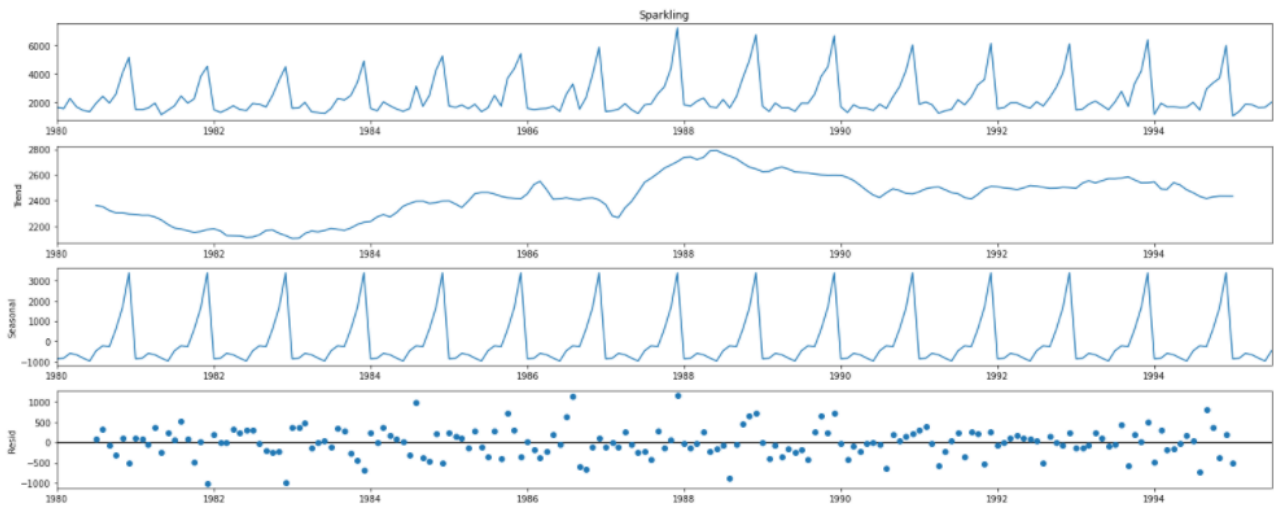


Fig.8

- Multiplicative Model

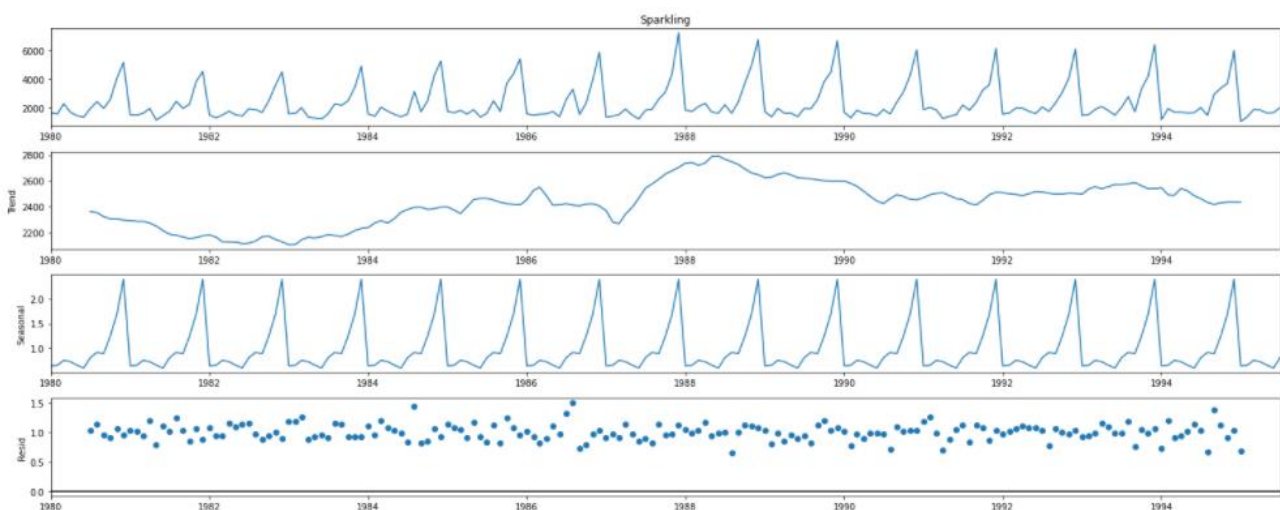
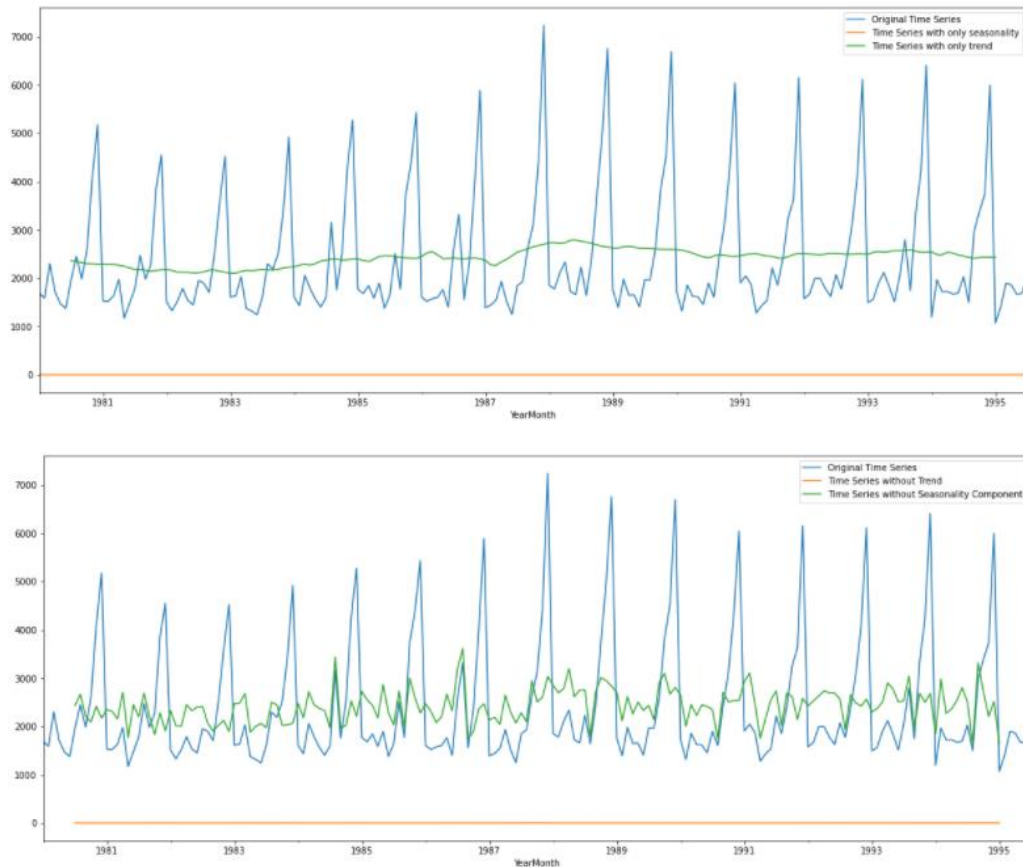


Fig.9

- The residual shows high variability across the period of time series, which is more or less consistent in both additive and multiplicative decompositions.
  - Residual mean of additive decomposition is : -1.2088
  - Residual mean of multiplicative decomposition is : 0.9974
- Out of the two multiplicative model is slightly better. P-value of shapiro test for multiplicative model is also significant. Our time series can be treated as multiplicative for model building.
- For the multiplicative series, we see that a lot of residuals are located around 1.



**Fig.10**

- The above graph shows the original time series, time series with only Seasonality and time series with only trend for multiplicative model.
- Second graph shows time series without trend and time series without seasonality.

**1.3** Split the data into training and test. The test data should start in 1991. Read the data as an appropriate Time Series data and plot the data.

- The time series is split into train and test set with test set starting from year 1991.
- Train Set

Sparkling	
YearMonth	
1980-01-31	1686
1980-02-29	1591
1980-03-31	2304
1980-04-30	1712
1980-05-31	1471

Train Head

Sparkling	
YearMonth	
1990-08-31	1605
1990-09-30	2424
1990-10-31	3116
1990-11-30	4286
1990-12-31	6047

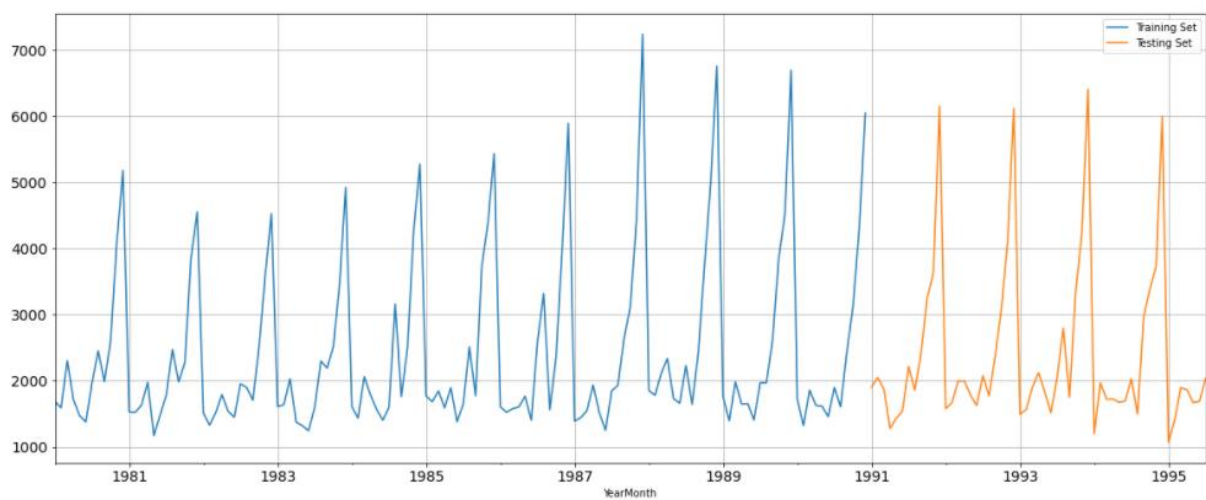
Train Tail

- Test Set

Sparkling		Sparkling	
YearMonth		YearMonth	
1991-01-31	1902	1995-03-31	1897
1991-02-28	2049	1995-04-30	1862
1991-03-31	1874	1995-05-31	1670
1991-04-30	1279	1995-06-30	1688
1991-05-31	1432	1995-07-31	2031

Test Head

Test Tail



**Fig.11**

- Plot of train and test set for the sparkling time series.

#### 1.4 Build various exponential smoothing models on the training data and evaluate the model using RMSE on the test data.

Other models such as regression, naïve forecast models, simple average models etc. should also be built on the training data and check the performance on the test data using RMSE. Read the data as an appropriate Time Series data and plot the data.

##### Linear Regression on Time:

- To regress the sale of sparkling wine, numerical time instance needs to be added to train and test set respectively for regression algorithm to work on our time series.
- Model is trained on training set and RMSE evaluated on test set.
  - ❖ Test RMSE – 1389.135

##### Naïve Forecasting:

- In naïve model, the value at the end of train set is applied to all the test set records. Prediction is made using these values with actual values in test set.
  - ❖ Test RMSE – 3864.279
- Model is very poor fitted and does not capture neither trend nor seasonality present in the dataset.

##### Simple Average Forecasting:

- In simple average model, the forecast is done using the mean of the target variable of training set of the time series.
  - ❖ Test RMSE – 1275.082
- Model is very poor fitted and does not capture neither trend nor seasonality present in the dataset.

##### Trailing Moving Average:

- For this model we will calculate the rolling means for different time intervals. The best interval can be determined by the least RMSE value.
- Trailing moving average models are built for 2, 4, 6 and 9 points moving averages.
  - ❖ Test RMSE - 2 Trailing Average – 813.401
  - ❖ Test RMSE - 4 Trailing Average – 1156.590
  - ❖ Test RMSE - 6 Trailing Average – 1283.927
  - ❖ Test RMSE - 9 Trailing Average – 1346.278
- The best trailing interval for moving average is 2.

##### Simple Exponential Smoothing-Auto fit:

- Simple exponential smoothing is applied if the time series has neither trend nor seasonality, which is not the case for our dataset.
- For SES auto fit we are just passing the train set to the base model and evaluating on test set.
  - ❖ Test RMSE – 1316.035

### Simple Exponential Smoothing Manual:

- Here we are manually finding the best value for alpha between 0 and 1.
- For different value of alpha, we are fitting it to the SES model and we are evaluating the test RMSE value. The best test RMSE value found is for 0.1 alpha.
  - ❖ Test RMSE – 1375.39

### Double Exponential Smoothing-Manual:

- Double exponential smoothing model is applied when the data has trend, but no seasonality, which is not the case for our dataset. Our dataset has slight trend and seasonality both.
- For DES model we trying to find the best value for both alpha and beta respectively. The best combination of both alpha and beta are chosen based on the test RMSE value. We get the best test RMSE value at alpha = 0.1 and beta = 0.1.
  - ❖ Test RMSE – 1777.73

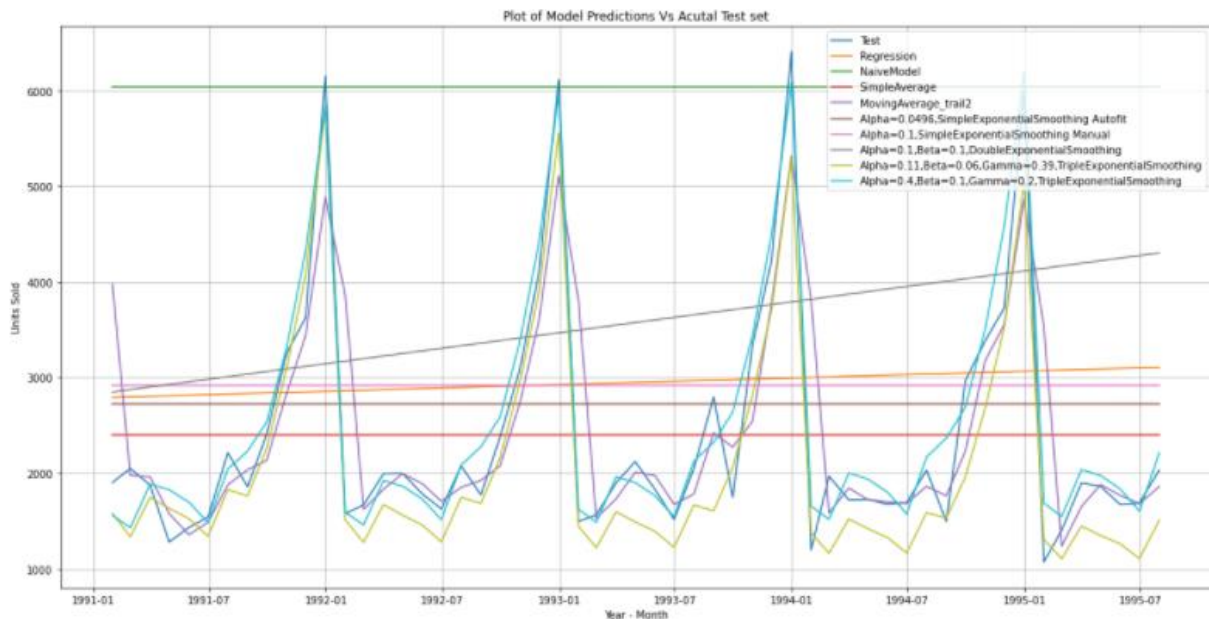
### Triple Exponential Smoothing-Auto fit:

- Triple exponential smoothing model is applied when data has both trend and seasonality. Our sparkling time series has both trend and seasonality present.
- For auto fit TES model we are just passing the train set to the base model and evaluating on the test set.
- We get the best value for alpha = 0.11, beta=0.06, gamma=0.39 from the auto fit model.
  - ❖ Test RMSE – 469.76

### Triple Exponential Smoothing Manual:

- For TES manual model we try to find the best value of alpha, beta and gamma from range between 0.1 -1 for each respectively. Trying to find the best combination of value and testing on test set to find the least RMSE value.
- We see that we get the least value for alpha=0.4, beta=0.1 and gamma=0.2.
  - ❖ Test RMSE – 336.71

## Model Comparison:



**Fig.12**

- We have plotted the test data against predictions from all the models.
- Out of all the models we can see that the best model is TES manual model. Test RMSE value for the TES manual model is also least among all the models built so far.

**1.5** Check for the stationarity of the data on which the model is being built on using appropriate statistical tests and also mention the hypothesis for the statistical test. If the data is found to be non-stationary, take appropriate steps to make it stationary. Check the new data for stationarity and comment. Read the data as an appropriate Time Series data and plot the data.

- Augmented Dickey Fuller test is the statistical test to check the stationarity of a time series. The test determine the presence of unit root in the series to understand if the series is stationary or not.
- For checking the stationary of the series we have the following hypothesis as below:  
H0: The series has a unit root, hence it is not stationary.  
H1: The series has no unit root, hence it is stationary.
- Now to test our hypothesis we perform the Dickey Fuller Test on the whole data.

**Output:**

DF test statistic is -1.798

DF test p-value is 0.7055958459932753

Number of lags used 12

Since p-value is significant (greater than 0.05) we fail to reject the null hypothesis. So data is not stationary.

- Taking differencing of order 1 on the whole series and again testing for stationarity.

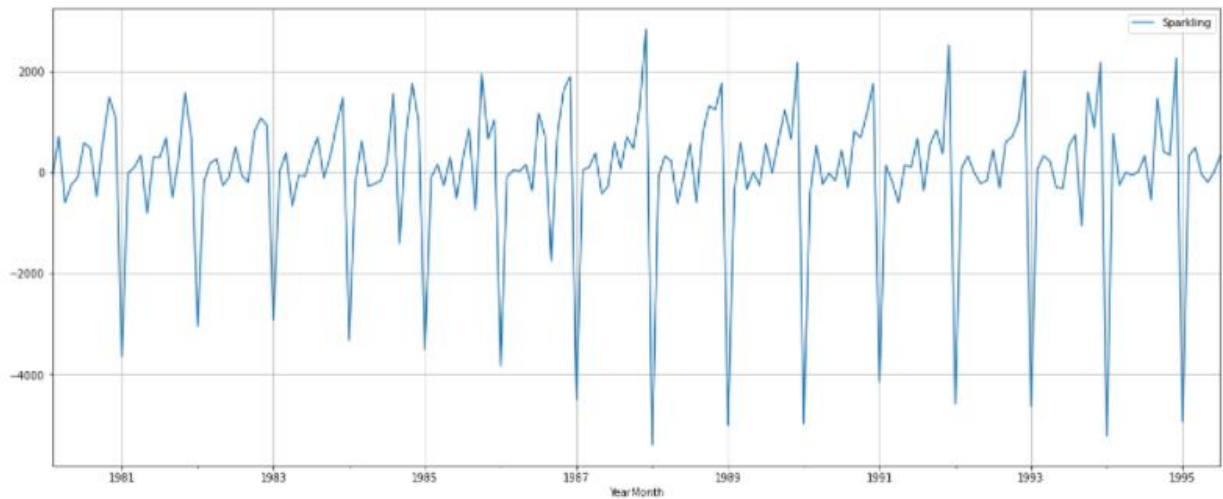
**Output:**

DF test statistic is -44.912

DF test p-value is 0.0

Number of lags used 10

Since the p-value is low less than  $\alpha=0.05$  we reject the null hypothesis. So for differencing of order 1 we get a stationary time series.



**Fig.13**

- We can see that the data has trend and seasonality component present in it but the series is now stationary around 0.
- We also check for stationarity of the train set and get the value of  $d=1$ .



- 1.6** Build an automated version of the ARIMA/SARIMA model in which the parameters are selected using the lowest Akaike Information Criteria (AIC) on the training data and evaluate this model on the test data using RMSE.

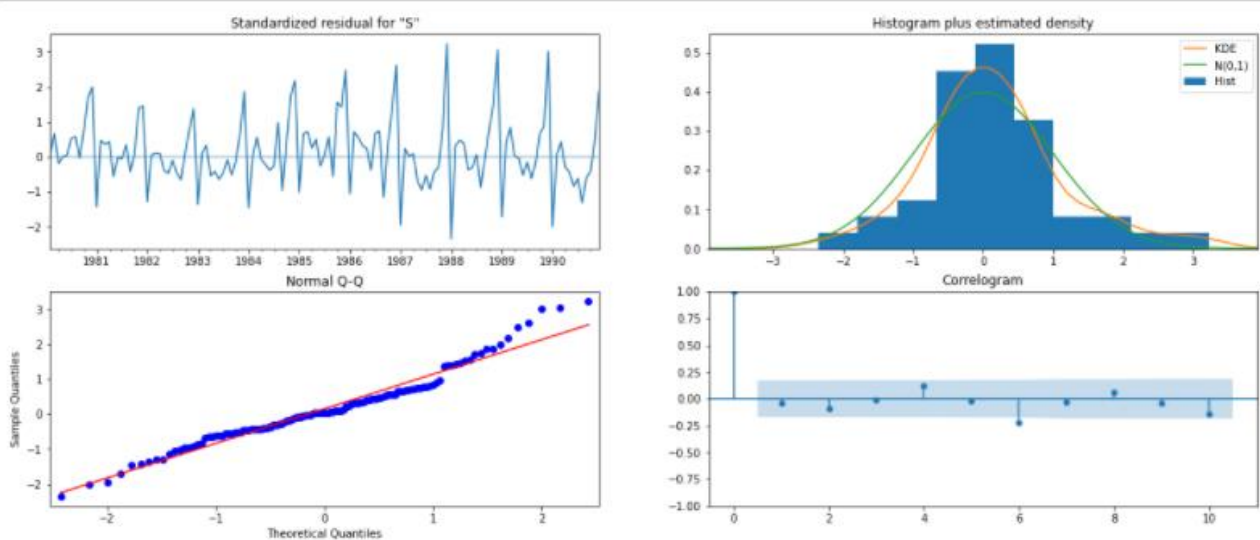
#### ARIMA:

- We build an automated ARIMA model on the dataset taking order of differencing as 1 and p and q value from range 0-4. Finding all combination of values of (p,d,q) and fitting in the ARIMA model to find the least AIC score.
- We get least AIC value for ARIMA(2,1,2).
  - ❖ AIC value for ARIMA(2,1,2) – 2213.50

SARIMAX Results						
=====						
Dep. Variable:	Sparkling	No. Observations:	132			
Model:	ARIMA(2, 1, 2)	Log Likelihood	-1101.755			
Date:	Wed, 12 Jan 2022	AIC	2213.509			
Time:	13:53:03	BIC	2227.885			
Sample:	01-31-1980	HQIC	2219.351			
	- 12-31-1990					
Covariance Type:	opg					
=====						
	coef	std err	z	P> z	[0.025	0.975]
-----						
ar.L1	1.3121	0.046	28.781	0.000	1.223	1.401
ar.L2	-0.5593	0.072	-7.741	0.000	-0.701	-0.418
ma.L1	-1.9917	0.109	-18.217	0.000	-2.206	-1.777
ma.L2	0.9999	0.110	9.109	0.000	0.785	1.215
sigma2	1.099e+06	1.99e-07	5.51e+12	0.000	1.1e+06	1.1e+06
=====						
Ljung-Box (L1) (Q):	0.19	Jarque-Bera (JB):	14.46			
Prob(Q):	0.67	Prob(JB):	0.00			
Heteroskedasticity (H):	2.43	Skew:	0.61			
Prob(H) (two-sided):	0.00	Kurtosis:	4.08			
=====						
Warnings:						
[1] Covariance matrix calculated using the outer product of gradients (complex-step).						
[2] Covariance matrix is singular or near-singular, with condition number 3.24e+27. Standard errors may be unstable.						

**Table.3**

- From the model summary it can be inferred that all AR and MA terms are significant, as their value are below 0.05.



**Fig.14 Diagnostics plot – ARIMA(2,1,2)**

- The diagnostics plot of model was derived and the standardized residuals are found to follow a mean of zero, and the histogram shows the residuals follow a normal distribution.
- Test RMSE for ARIMA(2,1,2) – 1299.97

#### SARIMA:

- We build an automated SARIMA model on the dataset taking order of differencing as 1 and P,p,Q,q value from range 0-4. We have taking seasonality value as 12, as data is of monthly period in the entire series and D=0. Finding all combination of values of (p,d,q)(P,D,Q,S) and fitting in the SARIMA model to find the least AIC score.
- The value of (p,d,q)(P,D,Q,S) for which SARIMA model is stable is (3,1,2)(3,0,1,12).
- There are other less AIC values as well but the residual histogram for those values is not normally distributed, so we are not considering those values.
  - ❖ AIC value for SARIMA (3,1,2)(3,0,1,12) – 1388.60

```

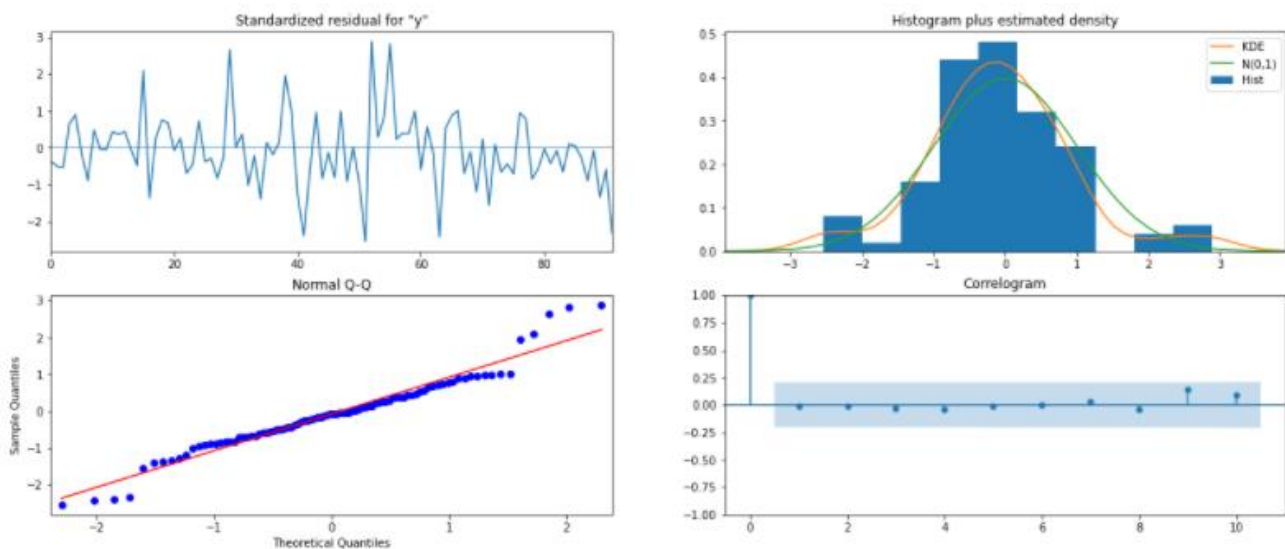
=====
SARIMAX Results
=====
Dep. Variable:          y          No. Observations:      132
Model:                 SARIMAX(3, 1, 2)x(3, 0, [1], 12)    Log Likelihood        -684.301
Date:                  Wed, 12 Jan 2022                  AIC                  1388.603
Time:                  13:59:38                          BIC                  1413.821
Sample:                0                                HQIC                 1398.781
Covariance Type:      opg
=====
              coef      std err      z      P>|z|      [0.025      0.975]
-----
ar.L1         -0.5434      0.416     -1.307     0.191     -1.358      0.271
ar.L2         -0.0078      0.198     -0.039     0.969     -0.396      0.381
ar.L3          0.0635      0.140      0.452     0.651     -0.212      0.339
ma.L1         -0.1992      0.404     -0.493     0.622     -0.991      0.593
ma.L2         -0.6547      0.326     -2.006     0.045     -1.295     -0.015
ar.S.L12       0.7653      0.448      1.707     0.088     -0.113      1.644
ar.S.L24       0.1091      0.330      0.331     0.741     -0.537      0.756
ar.S.L36       0.1763      0.186      0.946     0.344     -0.189      0.542
ma.S.L12      -0.2428      0.451     -0.539     0.590     -1.126      0.640
sigma2        1.663e+05    2.63e+04     6.328     0.000     1.15e+05    2.18e+05
=====
Ljung-Box (L1) (Q):      0.00    Jarque-Bera (JB):      9.36
Prob(Q):                0.96    Prob(JB):           0.01
Heteroskedasticity (H):  1.25    Skew:              0.35
Prob(H) (two-sided):    0.54    Kurtosis:          4.40
=====

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).

```

**Table.4**

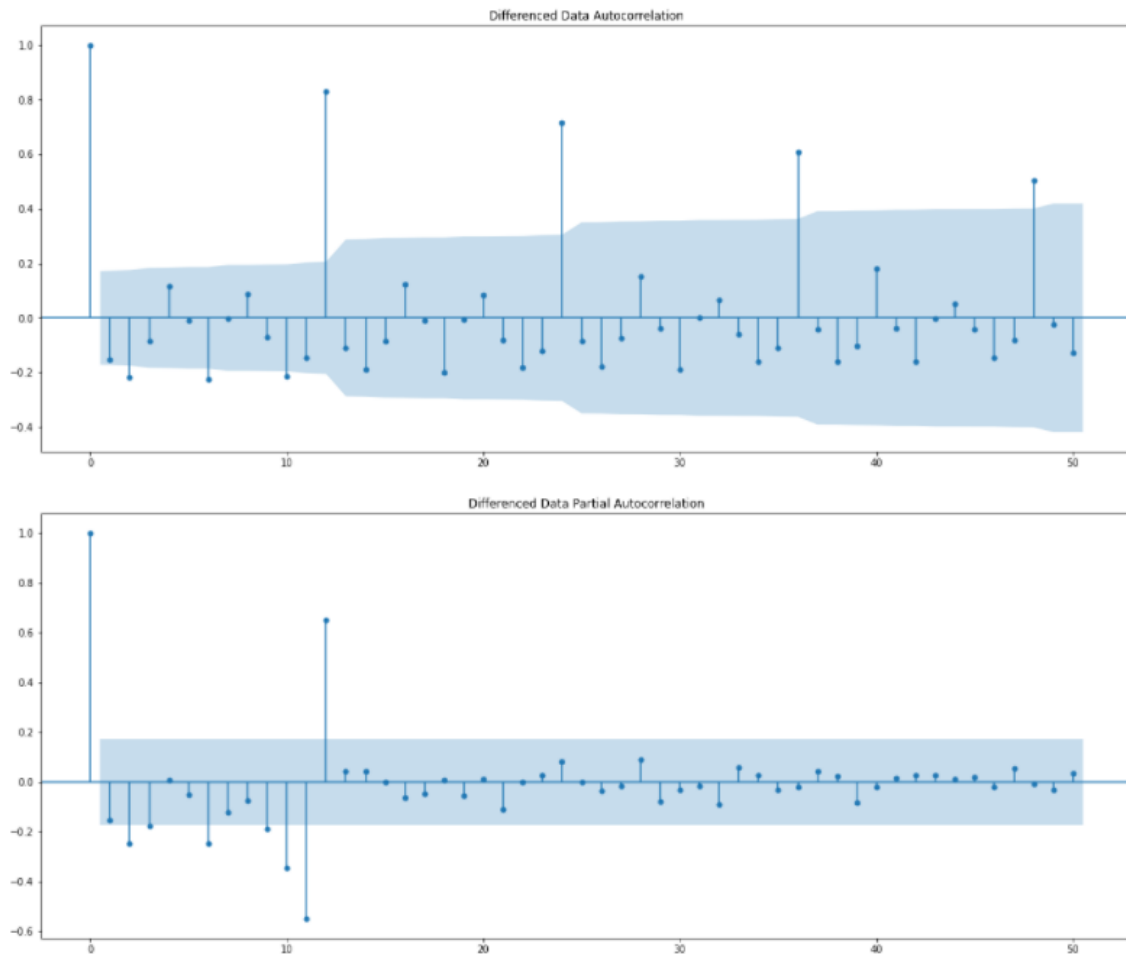
- From the model summary it can be inferred that most AR and MA terms are not significant.



**Fig.15 Diagnostics plot – SARIMA(3,1,2)(3,0,1,12)**

- Inference from model diagnostics confirms that the model residuals are normally distributed.
- Standardised residual** - Do not display any obvious seasonality.
- Histogram plus estimated density** –The KDE plot of the residuals is similar with the normal distribution, hence the model residuals are normally distributed.
- Normal Q-Q plot** - There is an ordered distribution of residuals(blue dots) following the linear trend of the samples taken from the standard normal distribution with  $N(0,1)$
- Correlogram** - The time series residuals have low correlation with lagged versions of itself.
- Test RMSE for SARIMA(3,1,2)(3,0,1,12) – 580.18

**1.7** Build ARIMA/SARIMA models based on the cut-off points of ACF and PACF on the training data and evaluate this model on the test data using RMSE.



**Fig.16**

- While reading the PACF and ACF plot:  
We always look at the positive side.  
We exclude 1 lag as it represents the series itself.  
Look at consecutive bars that exceed the threshold to find the values from graph.

**ARIMA Model Using PACF and ACF:**

- Best parameters are selected by looking at the ACF and the PACF plot
- Here, we have taken  $\alpha=0.05$ .  
The Auto-Regressive parameter in an ARIMA model is 'p' which comes from the significant lag before which the PACF plot cuts-off to 0.  
The Moving-Average parameter in an ARIMA model is 'q' which comes from the significant lag before the ACF plot cuts-off to 0.  
By looking at the above plots, we can say that both the PACF and ACF plot cuts-off at lag 0.  
Difference order is of 1.

- ARIMA Model using  $p, q = 0$  and  $d=1$ .

```

=====
SARIMAX Results
=====
Dep. Variable:      Sparkling    No. Observations:      132
Model:             ARIMA(0, 1, 0)  Log Likelihood         -1132.832
Date:              Wed, 12 Jan 2022  AIC                      2267.663
Time:              13:59:39         BIC                      2270.538
Sample:            01-31-1980       HQIC                     2268.831
                  - 12-31-1990
Covariance Type:    opg
=====
              coef      std err      z      P>|z|      [0.025      0.975]
-----
sigma2      1.885e+06  1.29e+05  14.658    0.000    1.63e+06  2.14e+06
=====
Ljung-Box (L1) (Q):                3.07  Jarque-Bera (JB):                198.83
Prob(Q):                           0.08  Prob(JB):                       0.00
Heteroskedasticity (H):             2.46  Skew:                          -1.92
Prob(H) (two-sided):               0.00  Kurtosis:                       7.65
=====

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).

```

Table.5

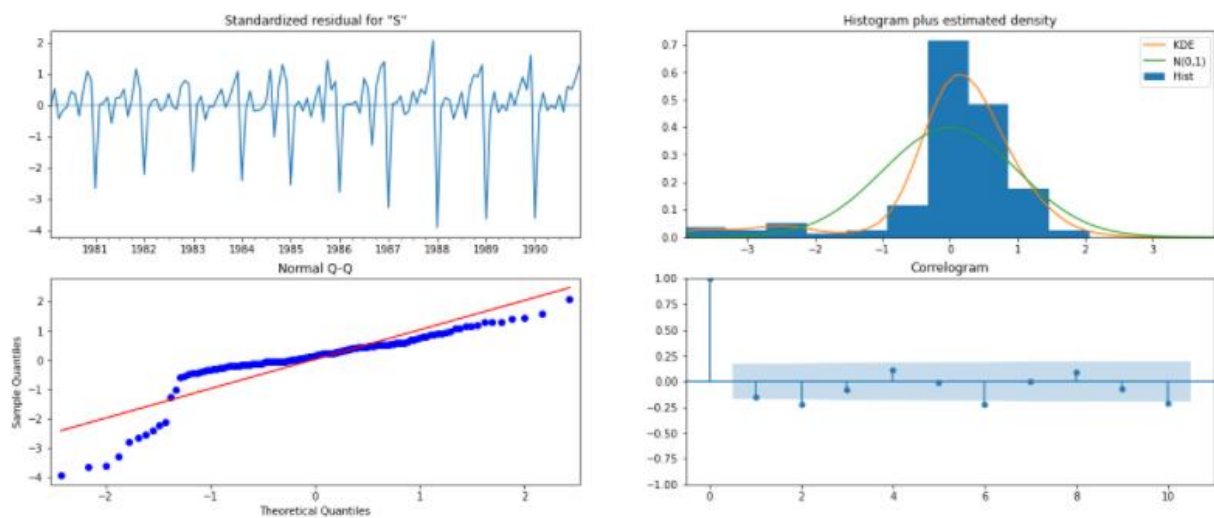


Fig.17

- Since we know that our time series has both trend and seasonality components ARIMA model thus built using  $p$  and  $q = 0$  is not correct at all.
- Test RMSE for ARIMA(0,1,0) – 3864.279
- As we can see from the test RMSE also the model is performing very poorly.

## SARIMA Model Using PACF and ACF:

- Here, we have taken  $\alpha=0.05$ .

We are going to take the seasonal period as 12. We will keep the  $p = 0$ ,  $q = 0$  and  $d=1$  parameters same as the ARIMA model.

The Auto-Regressive parameter in an SARIMA model is 'P' which comes from the significant lag after which the PACF plot cuts-off to 0.

The Moving-Average parameter in an SARIMA model is 'Q' which comes from the significant lag after which the ACF plot cuts-off to 0.

We have checked the ACF and the PACF plots only at multiples of 12 (since 12 is the seasonal period).

- SARIMA Model using  $p,d,q = 0,1,0$  and  $P,D,Q = 0,0,0,12$ .

```

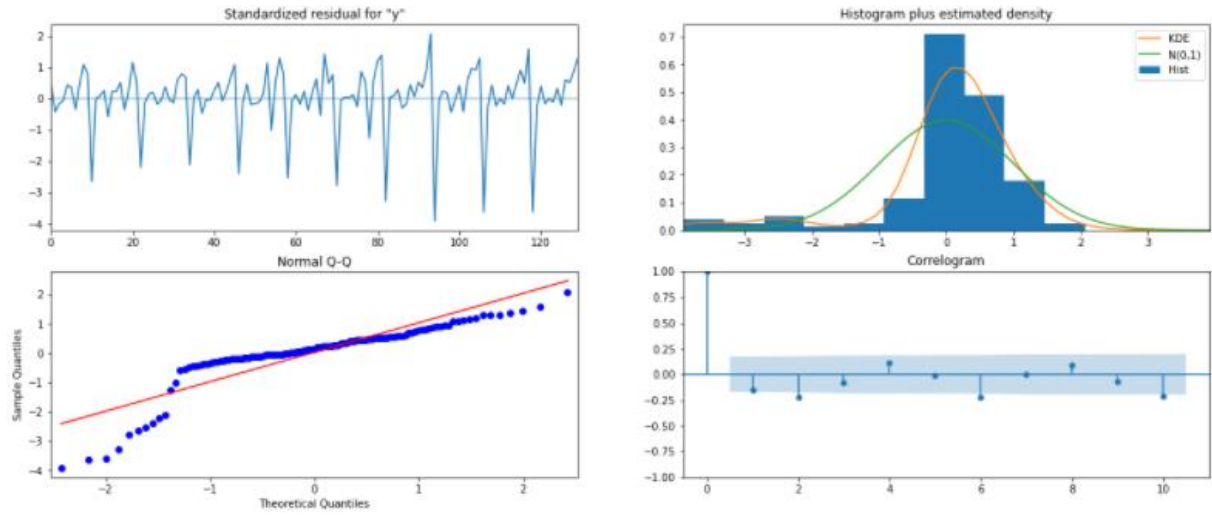
=====
SARIMAX Results
=====
Dep. Variable:          y      No. Observations:          132
Model:                SARIMAX(0, 1, 0)      Log Likelihood          -1124.680
Date:                Wed, 12 Jan 2022      AIC                  2251.360
Time:                13:59:40      BIC                  2254.227
Sample:              0      HQIC                  2252.525
                   - 132
Covariance Type:      opg
=====
              coef      std err          z      P>|z|      [0.025      0.975]
-----
sigma2      1.899e+06   1.31e+05    14.543     0.000     1.64e+06   2.16e+06
=====
Ljung-Box (L1) (Q):                3.04      Jarque-Bera (JB):                194.29
Prob(Q):                          0.08      Prob(JB):                  0.00
Heteroskedasticity (H):            2.46      Skew:                      -1.92
Prob(H) (two-sided):              0.00      Kurtosis:                   7.60
=====

```

Warnings:

[1] Covariance matrix calculated using the outer product of gradients (complex-step).

**Table.6**



**Fig.18**

- Again as the value of  $p, q$  and  $P, Q$  has been taken from the ACF and PACF plots our model is not generalising well and performing very poorly for the SARIMA model as well.
- Test RMSE for SARIMA(0,1,0)(0,0,0,12) – 3864.279

- 1.8** Build a table with all the models built along with their corresponding parameters and the respective RMSE values on the test data.

	Test RMSE
Alpha=0.4,Beta=0.1,Gamma=0.2,TripleExponentialSmoothing	336.715250
Alpha=0.11,Beta=0.06,Gamma=0.39,TripleExponentialSmoothing-Autofit	469.767970
pdq=(3,1,1),PDQS=(3,0,0,12),SARIMA Automated	580.189783
2pointTrailingMovingAverage	813.400684
4pointTrailingMovingAverage	1156.589694
Simple Average	1275.081804
6pointTrailingMovingAverage	1283.927428
Alpha=2,Beta=1,Gamma=2,ARIMA Automated	1299.979569
Alpha=0.0496,SimpleExponentialSmoothing-Autofit	1316.035487
9pointTrailingMovingAverage	1346.278315
Alpha=0.1,SimpleExponentialSmoothing	1375.393398
RegressionOnTime	1389.135175
Alpha=0.1,Beta=0.1,DoubleExponentialSmoothing	1777.734773
Naive Model	3864.279352
Alpha=0,Beta=1,Gamma=0,ARIMA Manual	3864.279352
pdq=(0,1,0),PDQS(0,0,0,12),SARIMA Manual	3864.279352

**Table.7**

- The above table provides test RMSE value for all the models that we have performed so far.
- We have arranged the table in ascending order of Test RMSE values.
- Model with lowest Test RMSE value will help in predicting the future sales of sparkling wine much better.
- The best model as per the table is Triple Exponential Smoothing manual model. Having the least RMSE value of 336.71.



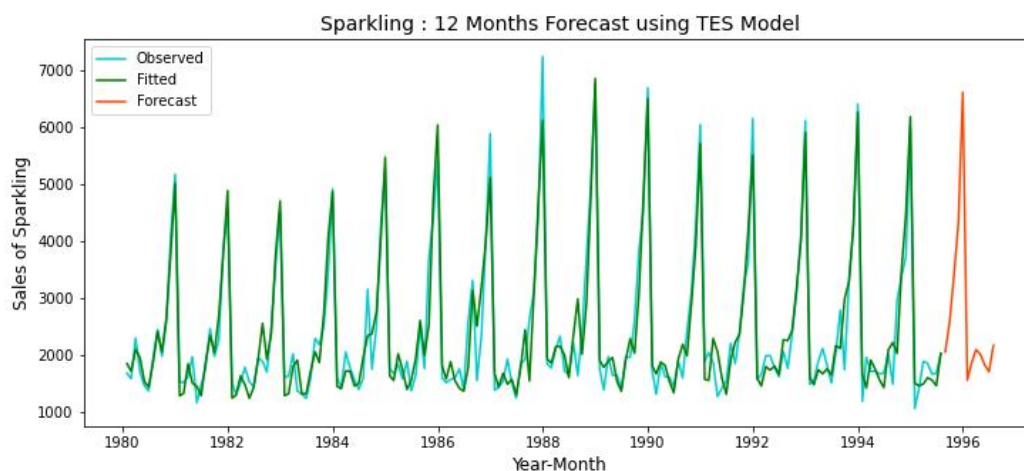
**1.9** Based on the model-building exercise, build the most optimum model(s) on the complete data and predict 12 months into the future with appropriate confidence intervals/bands.

- Now we built our best model that is Triple Exponential Smoothing on the complete dataset using the best parameters values we found earlier.
- Parameter values used for Triple Exponential Smoothing are:  
Alpha: 0.4  
Beta: 0.1  
Gamma: 0.2  
Optimized: True  
Use\_brute: True
- Now using the above model we predict for 12 month sales figure of 'Sparkling' wine in the future.

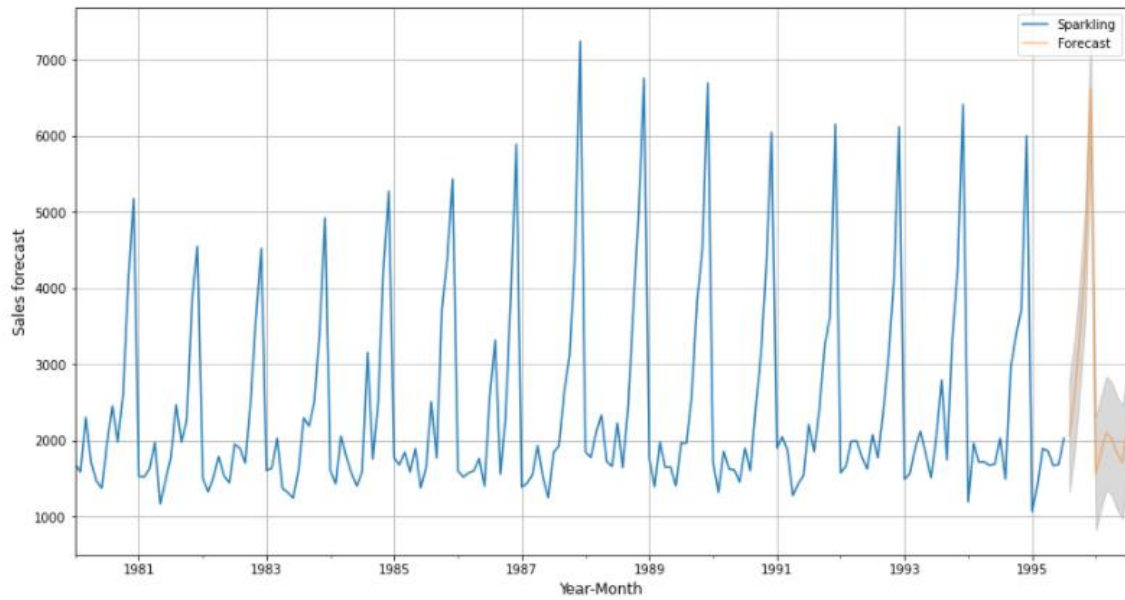
```
1995-08-31    2063.370030
1995-09-30    2579.777769
1995-10-31    3418.086343
1995-11-30    4308.589385
1995-12-31    6615.784148
1996-01-31    1566.940526
1996-02-29    1852.450225
1996-03-31    2101.024170
1996-04-30    2024.178682
1996-05-31    1835.548594
1996-06-30    1713.199103
1996-07-31    2177.175089
Freq: M, dtype: float64
```

**Table.8**

- Plotting our time series against the values found out by our model for the complete data and our prediction for 12 months in the future.
- We can visually see that our model is able to mimic the time series pattern quite reasonably well.

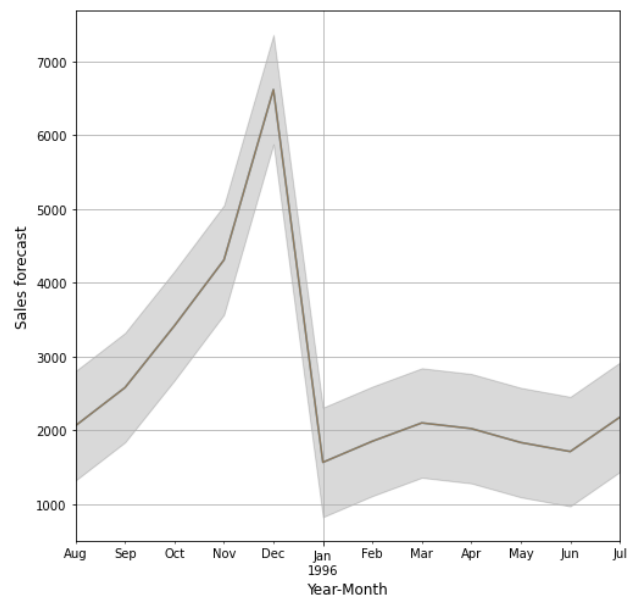


**Fig.19**



**Fig.20**

- Plot original data and of our prediction for 12 month sale in the future against confidence interval of 95%.



**Fig.21**

- Plot of our forecast along with the confidence interval band of 95%.

**1.10** Comment on the model thus built and report your findings and suggest the measures that the company should be taking for future sales.

- The model forecast sale of 32256 units of sparkling wine in 12 months in the future with average sale of 2688 units per month.
- Highest sales is predicted to occur in Dec-1995 with unit's sales hitting maximum of 6616 units.
- Sparkling sales seems to have stabilized over the years and not much trend compared to previous years.
- December month shows the highest Sales across the years from 1980-1995.
- The models are built considering the Trend and Seasonality into account and we see from the output plot that the future prediction is in line with the trend and seasonality in the previous years.
- The Sales of Sparling wine is seasonal, hence the company cannot have the same stock through the year. The predictions would help here to plan the Stock need basis the forecasted sales.
- The company should use the prediction results and capitalize on the high demand seasons and ensure to source and supply the high demand.
- The company should use the prediction results to plan the low demand seasons to stock as per the demand.
- The forecast also indicates that year-on-year the sales of sparkling wine is not showing an upward trend. The wine company must adopt innovative marketing skills to improve the sales.

**THE END**